

---

# *LiRo*: Benchmark and leaderboard for Romanian language tasks

---

<b>Stefan Dumitrescu</b>	<b>Petru Rebeja</b> AI Cuza University of Iasi	<b>Beata Lorincz</b> Technical University of Cluj-Napoca	
<b>Mihaela Gaman</b> University of Bucharest	<b>Andrei-Marius Avram</b> Politehnica University of Bucharest	<b>Mihai Ilie</b>	
<b>Andrei Pruteanu</b>	<b>Adriana Stan</b> Technical University of Cluj-Napoca	<b>Lorena Rosia</b> Deloitte	
<b>Cristina Iacobescu</b> Deloitte	<b>Luciana Morogan</b> Military Technical Academy	<b>George-Andrei Dima</b> Politehnica University of Bucharest	
<b>Gabriel Marchidan</b> Feel IT Services	<b>Traian Rebedea</b> Politehnica University of Bucharest	<b>Madalina Chitez</b> West University of Timisoara	
<b>Dani Yogatama</b> DeepMind	<b>Sebastian Ruder</b> DeepMind	<b>Radu Tudor Ionescu</b> University of Bucharest	<b>Razvan Pascanu</b> DeepMind
<b>Viorica Patraucean</b> DeepMind viorica@google.com			

## Abstract

1 Recent advances in NLP have been sustained by the availability of large amounts of  
2 data and standardized benchmarks, which are not available for many languages. As  
3 a small step towards addressing this, we propose *LiRo*, a platform for benchmarking  
4 models on the Romanian language on nine standard tasks: text classification,  
5 named entity recognition, machine translation, sentiment analysis, POS tagging,  
6 dependency parsing, language modelling, question-answering, and semantic textual  
7 similarity. We also include a less standard task of Romanian embeddings debiasing,  
8 to address the growing concerns related to gender bias in language models. The  
9 platform exposes per-task leaderboards populated with baseline results for each  
10 task. In addition, we create three new datasets: one from Romanian Wikipedia  
11 and two by translating the Semantic Textual Similarity (STS) benchmark and  
12 the Cross-lingual Question Answering Dataset (XQuAD) into Romanian. We  
13 believe *LiRo* will not only add to the growing body of benchmarks covering various  
14 languages, but can also enable multi-lingual research by augmenting parallel  
15 corpora, and hence is of interest for the wider NLP community. *LiRo* is available at  
16 <https://lirobenchmark.github.io/>

## 17 1 Introduction

18 Recent years have seen rapid progress on many language understanding tasks, from language mod-  
19 elling [e.g. 4] to translation [e.g. 27] or Q&A [e.g. 21]. Most of these understandably have happened

20 in English, relying on the proliferation of datasets [e.g. 7, 33] and on easy access to leaderboards  
21 and benchmarks<sup>1</sup> [e.g. 43] that facilitate communication and standardization of experiments. Unfor-  
22 tunately, a similar level of access is lacking for many other languages. In this work, we focus on  
23 Romanian and aim to provide datasets and tools to facilitate research on Romanian language tasks.

24 Romanian is an Indo-European Romance language that evolved in relative isolation compared to other  
25 Romance languages, leading to its unique characteristics. In particular, Romanian has mixed linguistic  
26 typology [8], displaying characteristics from two different families: Romance languages [23] and  
27 Balkansprachbund [39]. For example, the majority of verb forms in Romanian function syntactically  
28 as in other languages included in the Italian branch of the Indo-European Romance language family,  
29 with the shift from Latin to Romance manifesting as the shift from synthetic/inflectional towards  
30 analytic/syntagmatic constructions (e.g. Latin *feci*, Italian *ho fatto*, Romanian *am făcut*). However,  
31 the geographical proximity to the Balkan region accounts for the existence of verb forms, such as  
32 the volo future [26], that are common to Romanian and Slavic languages (e.g. Romanian *voi face*,  
33 Bulgarian *shte napravya*). Similarly, other features, such as the enclitic definite article, attached in  
34 Romanian at the end of the noun (e.g. *omul* → *the man*) can either be explained through post-Roman  
35 regional contact in the Balkans or the influence of the Ancient Greek on Vulgar Latin [12]. The  
36 case of double negatives (e.g. *nu am mâncat nimic*), also present in French, Spanish and Italian,  
37 represents a challenge for ML algorithms trained on English language, where double negation are  
38 rather infrequent (e.g. *haven't eaten anything*). Furthermore, lexical similarity analyses emphasize  
39 the particularity of Romanian within the Romance group [18]. These are all arguments that support  
40 the latest findings in cross-lingual NLP studies stating that typological properties of languages impact  
41 allegedly *language-agnostic* models [20]. Hence, evaluating cross-lingual models on Romanian can  
42 contribute to shedding light on their predictive performance.

43 Although Romanian is spoken by around 25 million speakers, it is still considered a low-resourced  
44 language in terms of digital resources and NLP tools [40]. Within the European Language Grid [34],  
45 Romanian is listed with only 129 resources, tools and services, as opposed to English (2342), Spanish  
46 (658) or German (777).<sup>2</sup> To address this issue, we propose *LiRo* (**L**imba **R**omână = Romanian  
47 Language), the first benchmark and leaderboard targeting models for Romanian language tasks.  
48 Currently, it includes nine standard tasks (text classification, named entity recognition, machine  
49 translation, sentiment analysis, part-of-speech tagging, dependency parsing, language modelling,  
50 question-answering, semantic textual similarity) and one less standard task of gender-debiasing of  
51 language embeddings. We included this latter task to state the importance of studying language biases  
52 in ML models [19] and to encourage research in this direction also for the Romanian language.

53 Along with the platform, we introduce three new datasets: RO-STS (Romanian translation of the  
54 Semantic Textual Similarity dataset [6]), XQuAD-ro (the Romanian component of the XQuAD  
55 dataset [1]), and Wiki-ro (Romanian Wiki for language modelling evaluation). For part-of-speech  
56 tagging and dependency parsing, we rely on the Romanian version of UD-RRT [2], but we propose a  
57 cross-genre training-vs-testing split in order to measure the robustness of existing systems to stylistic  
58 changes – a relevant task for Romanian language, which tends to change its form across domains.

59 We provide baseline results for all the tasks either by extracting results from the literature for existing  
60 datasets or by creating new baselines for the newly-created datasets and the newly-created splits. We  
61 analyse the results of the new baselines and point to directions of improvement.

## 62 2 Related work

63 One of the first initiatives for the common evaluation of disjoint Natural Language Understanding  
64 (NLU) tasks was the General Language Understanding Evaluation [GLUE; 43] benchmark. Wang et al.  
65 [43] gathered nine tasks including question answering, sentiment analysis, and textual entailment, as  
66 well as their associated training and test datasets. GLUE also includes a diagnostic dataset to analyze  
67 models' performance with respect to a wide range of linguistic phenomena found in natural language.  
68 However, the rapid advancements in deep learning led to a quick saturation of the benchmark [42]  
69 where several models surpassed non-expert humans. Wang et al. [42] proposed SuperGLUE, a novel

---

<sup>1</sup>This includes websites such as [paperswithcode.com](https://paperswithcode.com).

<sup>2</sup>As of June 7, 2021, in the ELG Release 2, at <https://live.european-language-grid.eu/catalogue/>.

70 benchmark that includes a more diverse and challenging set of tasks. Additionally, SuperGLUE can  
71 showcase significant performance gaps between BERT-like models [13] and humans.

72 McCann et al. [29] introduced the Natural Language Decathlon (DecaNLP), a benchmark that  
73 comprises ten NLP tasks ranging from machine translation, question answering and summarization to  
74 sentiment analysis, relation extraction and semantic parsing. Poliak et al. [31] introduced the Diverse  
75 Natural Language Inference Collection (DNC), comprising 8 tasks and 13 existing datasets. DNC is  
76 aimed at evaluating a model’s capability to perform various types of reasoning. Another landmark  
77 collection of datasets for the English language was proposed by Conneau and Kiela [10]. SentEval  
78 [10] is advertised as a toolkit for the centralized evaluation of universal sentence representations. It  
79 is composed of 7 distinct tasks and 13 datasets. Different from the previous benchmarks, Evaluating  
80 Rationales And Simple English Reasoning (ERASER) [14] is a benchmark aiming to assess the  
81 interpretability of NLP models. The main contribution of this benchmark is the design of novel  
82 evaluation metrics to measure the alignment between human and model rationals. DeYoung et al.  
83 [14] establish that a rational is the evidence that supports a decision.

84 The aforementioned benchmarks are all based on English datasets. Recently, some effort has been  
85 dedicated to the development of multi-lingual benchmarks. XTREME [22] is a benchmark dedicated  
86 to the evaluation of cross-lingual generalization on 40 languages. Perhaps the most important  
87 observation of Hu et al. [22] is that state-of-the-art models for English exhibit sizeable performance  
88 gaps when transferred across languages. While the number of languages and the size of XTREME  
89 is remarkable, we emphasize that Romanian is not included. Through *LiRo*, we aim to establish a  
90 NLU benchmark for Romanian. Among the datasets included in the XTREME benchmark is the  
91 Cross-lingual Question Answering Dataset [XQuAD; 1] for which we provide a translation into  
92 Romanian by professional human translators, which was also added to the official XQuAD repository.

93 While some research works went towards creating multi-lingual benchmarks, other works focused on  
94 building mono-lingual benchmarks for understudied languages. For instance, a recently developed  
95 language-dependent benchmark is the Polish version of GLUE, known as the KLEJ benchmark [35].  
96 KLEJ contains a set of 9 evaluation tasks for the Polish language understanding. The authors collated  
97 existing datasets together with a new dataset for sentiment analysis. The platform provides evaluation  
98 code and a public leaderboard. Another example of mono-lingual NLU evaluation is IndoNLU [44],  
99 a benchmark dedicated to the Indonesian language. IndoNLU is composed of twelve tasks. The  
100 diversity of the tasks is ensured by selecting datasets from various domains and with different styles.  
101 Recently [30] proposed KLUE as a benchmark for the Korean language, also modelled after the  
102 GLUE benchmark.

103 Our platform currently includes 10 tasks, 8 datasets (out of which three are new) and a public  
104 leaderboard. We pledge to further develop *LiRo* and include additional datasets and tasks to provide a  
105 comprehensive evaluation platform for Romanian and multi-lingual language tasks.

### 106 3 *LiRo* benchmark and leaderboard

107 **Benchmark.** *LiRo* is an open-source benchmark and a continuous-submission leaderboard, concen-  
108 trating public Romanian datasets (existing and new) in specific tasks. The integration of datasets and  
109 tasks with model performance and efficiency allows both academia and industry to quickly gauge  
110 performance on tasks of interest. The benchmark also provides an overview of the Romanian NLU  
111 SoTA and direct access to relevant papers. Finally, it intends to foster a constructive competition and  
112 innovation by bringing together and promoting previously disparate resources.

113 *LiRo* is structured into *areas*, *tasks*, and *datasets*. In this paper, we focus on the NLP area, but in  
114 the future we intend to extend *LiRo* to other areas like speech or image captioning. Each area can  
115 have any number of tasks and for each task we can have any number of datasets, each with their  
116 own metric(s). *LiRo*’s homepage lists all available tasks, grouped by area. Each task contains a  
117 succinct description and the available datasets. A dataset is a specific corpus with defined training  
118 and evaluation splits, together with evaluation metrics and scripts to compute these metrics. A dataset  
119 can belong to multiple tasks—for example the Universal Dependencies Romanian RRT Treebank  
120 dataset [2] is used in POS tagging and parsing tasks. To keep things simple, *LiRo* does not host the  
121 datasets directly. Instead, we link to each individual resource’s webpage while having a dedicated  
122 description page for each dataset, with statistics about the dataset, metrics, and other details useful

#	Task	Dataset	Metrics	Score	Baseline
1.	Text Categorization by Topic	MOROCCO	Macro F1	88.03	[17]
2.	Named Entity Recognition	RONEC v1.0	Exact Match F1	85.88	[15]
3.	Machine Translation	WMT-16-ro-en	BLEU, ROUGE-L	38.5	[28]
4.	Sentiment Analysis	LaRoSeDa	F1	54.30	[38]
5.	POS Tagging	UD Ro-RRT (cross)	UPOS F1, XPOS F1	95.73	this paper
6.	Dependency Parsing	UD Ro-RRT (cross)	UAS F1, LAS F1	88.97	this paper
7.	Language Modelling	Wiki-ro	Perplexity	28.0	this paper
8.	Question Answering	XQuAD-ro	F1, EM	83.56	this paper
9.	Semantic Textual Similarity	RO-STs	Pearson, Spearman	0.81	this paper
10.	Gender debiasing	Ro embeddings	Modified-WEAT	2.57	this paper

Table 1: Tasks, datasets, associated metrics, and baseline results available in *LiRo*. Where there is more than one metric, only the result for the first one is reported here to reduce clutter. At the moment, *LiRo* contains 10 tasks with associated datasets. The baseline results for the first 4 tasks are from the top performing models existing in the literature (and included in *LiRo*), whereas for the remaining 6 tasks, we propose new datasets or new dataset splits and associated baselines.

123 for anyone who wishes to use them. For the newly-created datasets, we include details regarding  
124 licensing.

125 **Leaderboard.** Each dataset has its own leaderboard, both graphically displayed as an interactive  
126 chart, and as a table listing all participating models. For each model, we include (1) the rank of the  
127 model in the leaderboard, (2) model name, (3) metric values, (4) whether the model was trained on  
128 extra training data, (5) model size (number of parameters), (6) link to the model’s paper and online  
129 code repository if any, and (7) submission date. In contrast to other benchmarks, we decided to  
130 require model size as a first step towards evaluating not only performance but also computational  
131 efficiency, following recent trends focusing on green AI [37].

132 We chose to have a separate leaderboard per dataset. Other platforms formulate all tasks in a common  
133 setting (e.g. convert all tasks into a binary classification [35]), so that they can provide an aggregated  
134 score. However, we found that this can lead to artificial tasks and opaque scores that might not  
135 capture the performance of the models in a meaningful way, harming understanding. Hence, we  
136 decided to create separate leaderboards and use standard problem formulations and metrics.

137 To submit a model to the leaderboard, we provide a templated submission form that users have to  
138 fill in. The maintainers of the platform then request additional info if needed. Once a submission is  
139 approved by a maintainer, the new model’s results will be automatically displayed on the website. A  
140 similar process is used for submitting new tasks or datasets to the leaderboard.

### 141 3.1 Available Tasks

142 We list below the tasks currently included in the benchmark and their associated datasets and metrics.  
143 For a summary, see Table 1.

144 **1. Text Categorization by Topic:** is the task of assigning a sentence or document to an appropriate  
145 category. Currently, *LiRo* contains the MOROCCO dataset [5] with a Romanian and Moldavian news  
146 classification task.

147 **2. Named Entity Recognition:** is the task of identifying and labeling entities in a text with their  
148 corresponding type (e.g. person, date, location, etc.). We use RONEC [16], a fine-grained dataset of  
149 5,127 sentences annotated with 16 classes, totalling 26,376 annotated entities.

150 **3. Machine Translation:** is the task of translating a sentence from a source language to a different  
151 target language. Currently this task includes WMT16 RO-EN dataset [3], a classic translation corpus  
152 used in several NLP papers.

153 **4. Sentiment Analysis:** requires classifying the affective state of a text, most frequently labelled as  
154 positive or negative. We include the recently proposed LaRoSeDa dataset [38], the first and only  
155 public dataset to our knowledge for this task in Romanian.

156 **5. Part-of-Speech Tagging:** (POS tagging) is the task of tagging a word in a text with its part of  
157 speech. We use the standard Romanian dataset for this task, the Universal Dependencies Romanian  
158 RRT Treebank (UD-RRT) [2], but we propose a different train-test split of the data in order to evaluate

159 robustness across genres (see details in Section 5). UD-RRT has annotations for Universal Parts of  
160 Speech (UPOS) as well as language-specific parts of speech (XPOS).

161 **6. Dependency Parsing:** is the task of extracting a dependency parse of a sentence that represents its  
162 grammatical structure and defines the relationships between “head” words and words, which modify  
163 those heads. We use again UD-RRT with the same splits as in Task 5. UD-RRT offers multiple layers  
164 of annotation for dependency parsing.

165 **7. Language Modeling:** is the task of predicting the next word or character in a document. For  
166 evaluating models on this task, we release the Romanian Wiki dataset, described in the next section.

167 **8. Question-answering (QA):** The task is to answer a question given a segment of text as context. As  
168 the first such dataset in Romanian, we introduce XQuAD-ro, the Romanian translation of XQuAD [1].  
169 XQuAD follows the standard SQuAD [1, 33] setting for QA: given a context paragraph, the model  
170 has to answer questions whose answers (of variable length) are spans in the context paragraph.  
171 XQuAD-ro is further detailed in the next section.

172 **9. Semantic Textual Similarity:** Given a pair of sentences, this regression task measures how similar  
173 the sentences are. We introduce RO-STC as the Romanian translation of STC [6], see next section for  
174 more details on this dataset.

175 **10. Language embeddings debiasing:** Given the growing concern about the negative impact that  
176 gender-biased language embeddings may have in practical applications, we measure the gender bias  
177 in existing Romanian language embeddings using the method proposed in [45] for languages with  
178 grammatical gender, and invite contributors to submit debiasing methods that can lower the gender  
179 bias in existing embeddings, or submit less biased embeddings. More details in section 5.

## 180 4 Newly-proposed datasets

181 We introduce three new datasets: *RO-STC*, *XQuAD-ro*, and *Wiki-ro*. The *RO-STC* and *XQuAD-ro*  
182 datasets were carefully translated from English and are the first of their kind for Romanian. The  
183 *Wiki-ro* is the first officially published Wiki dump for the Romanian language, purposely-cleaned  
184 with the aim of standardizing language model evaluation.

185 **RO-STC.** The *RO-STC* (Romanian Semantic Textual Similarity) dataset is the Romanian translation  
186 of the STC English dataset<sup>3</sup>. RO-STC contains 8,628 sentence pairs with their similarity scores. The  
187 original English sentences were collected from news headlines, captions of images and user forums,  
188 and are categorized accordingly. The Romanian release follows this categorization and provides  
189 the same train/validation/test split with 5,749/1,500/1,379 sentence pairs in each subset. Using both  
190 translations and similarity scores, *RO-STC* can be used for (at least) two purposes: (1) as a textual  
191 similarity dataset for Romanian, and (2) as a parallel Romanian-English dataset that can be used in  
192 any downstream NLP task, e.g. machine translation. RO-STC contains 212,619 tokens out of which  
193 23,425 are unique. The average character length for the sentences is 66.39. The similarity scores  
194 for the sentences range from 0 to 5, with an average of 2.60. RO-STC is freely available in both the  
195 textual-similarity and the parallel corpus formats.<sup>4</sup>

196 To create the dataset, we first (i) obtained automatic translations using Google’s translation engine.  
197 Then, (ii) the data was partitioned, checked, and corrected by 10 volunteers (ML researchers for  
198 whom Romanian is their native language and speak English fluently). These corrected partitions  
199 were then (iii) assigned to 3 volunteers (Romanian linguistic master students) for final validation.  
200 The volunteers in both phases (ii) and (iii) received the original English sentences and the Romanian  
201 translations from the previous phase with the instruction: “correct the translation if needed to make it  
202 sound like natural Romanian whilst keeping the meaning as close as possible to the original English  
203 version”. We provide here BLEU scores to give an idea about the volume of modifications made in  
204 the two phases: BLEU(Google translations, final) = 62.8. BLEU(first correction, final) = 77.9. This  
205 shows that the initial automatic translation was very good, and the corrections made by volunteers  
206 improved the quality even more.

207 **XQuAD-ro.** The *XQuAD-ro* dataset contains the Romanian translations for the 240 paragraphs and  
208 1,190 question-answer pairs of the XQuAD [1] dataset, previously available for 11 languages. We

<sup>3</sup><https://ixa2.si.ehu.es/stswiki/index.php/STCbenchmark>

<sup>4</sup>Available at: <https://github.com/dumitrescustefan/RO-STC>.

Task	In-domain	Cross-domain	XQuAD-ro	F1	EM
POS Tagging	98.18	95.73	mBERT	72.69	58.99
Dependency Parsing	90.38	88.97	XLm-R Large	83.56	69.66

Table 2: Results for tasks on UD-RRT using the original in-domain splits and the proposed cross-domain splits.

Table 3: Zero-shot QA on XQuAD-ro.

209 obtained the Romanian version with the help of professional human translators. The average number  
 210 of tokens is 153.91 per paragraph, 12.03 per question and 3.33 per answer. The total number of  
 211 tokens is 55,229, with 10,570 unique tokens. The average number of questions per paragraph is 4.95.  
 212 The average character length of the paragraphs is 878.44, 67.01 for questions and 20.91 for answers.  
 213 *XQuAD-ro* is already included in the official XQuAD repository for free public access.

214 **Wiki-ro.** The *Wiki-ro* dataset contains the July 2020 dump of the Romanian Wikipedia. It was  
 215 thoroughly cleaned, with several custom rules. Besides removing all the wiki markup, we skipped  
 216 wiki pages that have a large quantity of sequential numbers—there are many documents that are  
 217 simply lists of years and events, unsuitable to calculate the perplexity of a language model. Other  
 218 rules include limiting foreign words, punctuation, very short documents, proxy documents, etc. The  
 219 corpus was segmented at the sentence level and tokenized, and is formatted as a one-sentence-per-line,  
 220 with empty lines delimiting documents. The dataset is divided into train, validation, and test splits,  
 221 always making sure that a document is entirely included in a single split. The train, validation, and  
 222 test sets have 2.1M lines and 44M words, 14K lines and 276K words, and 16K lines and 327K words,  
 223 respectively. The goal of this dataset is to provide standardized fine-tuning and evaluation of language  
 224 modelling of Romanian text.

## 225 5 Experiments

### 226 5.1 Cross-genre splits for UD-RRT

227 UD-RRT [2] contains texts from 9 different genres: Academic, FrameNet, Journalistic, Law, Litera-  
 228 ture, Medical, Miscellanea, Science, and Wikipedia. The original dataset contains *within domain*  
 229 train/valid/test splits for all 9 genres. Given the variability of the Romanian language across domains  
 230 (caused by the use of specific vocabulary terms and phrases), we propose to use this dataset in a  
 231 cross-domain setting for the tasks of dependency parsing and POS tagging, to better test the robust-  
 232 ness of language models. To this end, we consider Miscellanea as the test domain, while using the  
 233 remaining 8 domains for training. We chose Miscellanea as a test domain as it contains texts from all  
 234 the other domains, plus some extra domains, e.g. dictionary definitions. This makes Miscellanea a  
 235 good test set to probe generalisation, including to out-of-distribution samples. In Table 2, we can  
 236 observe that the baselines are not sufficiently robust across domains, losing about 2% in accuracy  
 237 on all the tasks. We used the Stanza framework [32] with default settings to run the in-domain and  
 238 cross-domain experiments.

### 239 5.2 Cross-lingual Q&A baseline

240 For the XQuAD-ro dataset, we provide the same baseline results as the original XQuAD paper [1].  
 241 Namely, we use mBERT [13] and XLM-R Large [9] trained on the English SQuAD v1.1 training  
 242 data and evaluate them via zero-shot transfer on XQuAD-ro test dataset. We report F1 and EM  
 243 (exact match) in Table 3. XLM-R Large significantly outperforms mBERT. This is not surprising  
 244 considering that the training set for XLM-R Large includes far more training data for Romanian  
 245 than mBERT’s training set: by volume of training data, Romanian is the 11th language for XLM-R  
 246 Large and the 30th language for mBERT. In fact, out of the 12 languages present in XQuAD, XLM-R  
 247 Large obtains the best results on Romanian, after English, in terms of both F1 and EM. The Russian  
 248 influences present in Romanian and the fact that Russian is the second language by volume in XLM-R  
 249 Large’s training set might explain this performance.

Model	Pearson coeff.	#params	Training	WMT16	RO-STS
RNN	0.6853	15M	RO-STS	2.9	21.9
ro-BERT (cased)	0.7927	124M	WMT16	24.7	30.9
ro-BERT (uncased)	<b>0.8159</b>	124M	WMT16 + RO-STS	<b>24.8</b>	44.0
mBERT (cased)	0.7664	167M	RO-STS Finetuned	24.6	<b>45.9</b>
mBERT (uncased)	0.7690	167M			

Table 4: RO-STS baselines for semantic similarity.

Table 5: Translation results on WMT16 and RO-STS test sets.

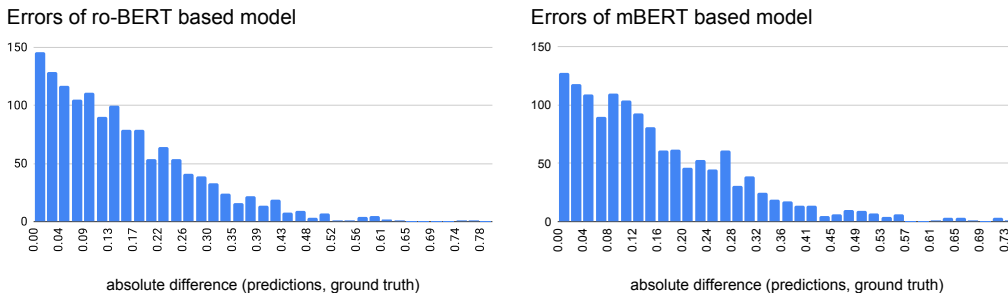


Figure 1: Errors made by two BERT-based models on the newly-created RO-STS dataset.

### 250 5.3 RO-STS baselines

251 For RO-STS dataset, we provide baselines for two tasks: Romanian semantic textual similarity and  
 252 EN  $\rightarrow$  RO translation, given the parallel nature of the dataset.

253 **Semantic textual similarity.** We include three semantic similarity baselines: an RNN-based model  
 254 and two transformer-based models, one using a monolingual Romanian BERT [ro-BERT; 15] and  
 255 one using a multilingual BERT [mBERT; 13]. The RNN-based model uses a two-layer bidirectional  
 256 LSTM to encode each sentence. Then, each sentence representation is passed through a standard  
 257 additive attention layer. For the transformer models, we encode each sentence separately, then  
 258 mean-pool the output token vectors. For all models, the similarity of the two resulting sentence  
 259 representations is computed using the cosine distance. This similarity is then compared with the  
 260 ground-truth scores normalized to  $[0, 1]$ . We use WordPiece tokenization and MSE loss for training.  
 261 For the BERT-based models, we experimented with both the ‘cased’ and ‘uncased’ datasets.

262 The results of the three models are included in Table 4, together with their size. The RNN model is  
 263 outperformed by the Transformer-based models in terms of Pearson coefficient. This is not surprising  
 264 given that the RNN-based model was trained from scratch and has a much lower capacity. Figure 1  
 265 shows histograms of the errors made by the two transformer-based models. Note that ro-BERT  
 266 is slightly more accurate than mBERT, the former having a more peaked histogram. In terms of  
 267 normalized absolute similarity error, ro-BERT obtained 0.154 and mBERT 0.160.

268 **Translation.** We provide a baseline for RO-STS as a parallel corpus. We employ WMT16 RO-EN  
 269 translation dataset [3] as a companion corpus and run the following experiments: (1) train on RO-STS,  
 270 (2) train on WMT16, (3) train on both WMT16 and RO-STS, and (4) train on WMT16 and finetune  
 271 on RO-STS. For modeling, we use the Open Neural Machine Translation (OpenNMT) toolkit [24]  
 272 and we employ the original Transformer model [41]. The sentences were encoded using the Unigram  
 273 subword tokenization [25], and we created a vocabulary of 8000 tokens for RO-STS training set and  
 274 a vocabulary of 32000 tokens for the rest. The sentences were batched together by their approximate  
 275 number of tokens resulting in batches of up to 2048 tokens for source and target sentences.

276 We evaluate the models from our four settings on WMT16 and RO-STS test sets and measure their  
 277 corresponding BLEU scores (see Table 5). The model trained on WMT16 obtains a BLEU score  
 278 of 24.7 on WMT16 test and a BLEU score of 30.9 on RO-STS test. On the other hand, the model  
 279 trained on RO-STS obtains a decent performance of 21.9 BLEU on RO-STS, but its performance is  
 280 dramatically reduced on WMT16 test due to the small size of the training dataset and vocabulary,  
 281 and domain mismatch. When training on both RO-STS and WMT16, the results on RO-STS were  
 282 significantly improved by 13.1 BLEU, while the results on WMT16 were just slightly improved with  
 283 0.1 BLEU. The highest BLEU score on RO-STS was achieved by the model that was first trained on

Sentence pair (En translation provided for reference)	Sim	ro-BERT	mBERT
Overestimating the similarity			
<i>(Un bărbat dansează, Un bărbat și o femeie dansează)</i> In English: <i>(A man dances, A man and a woman dance)</i>	0.4	0.76	0.80
<i>(Un pisoi bea lapte dintr-un bol, Un copil mic bea apă dintr-o cană)</i> <i>(A kitten drinks milk from a bowl, A small child drinks water from a cup)</i>	0.16	0.69	0.50
<i>(Nu ai nevoie de nicio viză, Nu ai nevoie de niciun fel de sos)</i> <i>(You don't need any visa, You don't need any kind of sauce)</i>	0	0.31	0.49
Underestimating the similarity			
<i>(Te-ai prins, Ai înțeles bine)</i> <i>(You got it, You understood well)</i>	1	0.40	0.31
<i>(Un bărbat râde cu o femeie, Un bărbat și o femeie râzând)</i> <i>(A man laughs with a woman, A man and a woman laughing)</i>	0.96	0.37	0.44
<i>(Ești pe drumul cel bun, Ai perfectă dreptate)</i> <i>(You are on the right track, You are perfectly right)</i>	0.8	0.23	0.07

Table 6: Example of errors made by the baseline models in predicting the similarity of sentences from RO-STS test set. The 2nd column ‘Sim’ is the ground truth, with 0 meaning no relation between the sentence pair and 1 meaning perfectly similar. 3rd and 4th columns are ro-BERT’s and mBERT’s predictions.

284 WMT16 and then finetuned on RO-STS, outperforming the previous model by 1.9 BLEU. However,  
285 as a result of fine-tuning on RO-STS, its performance slightly decreased on WMT16 by 0.1 BLEU.

#### 286 5.4 Wiki-ro baseline

287 We run zero-shot evaluation with a pre-trained ro-BERT masked language model [15], calculating  
288 pseudo-loglikelihood scores (PLLs) and their corresponding pseudo-perplexities (PPPLs) as in [36],  
289 obtaining: 29.08 (P)PPL on the validation set and 28.00 (P)PPL on the test set. This is a first, modest  
290 baseline which could be significantly improved, e.g. by fine-tuning the model on Wiki-ro training set.

#### 291 5.5 Gender debiasing baseline

292 We measure the gender bias in existing Romanian language embeddings [11] using the method  
293 proposed in [45] for languages with grammatical gender. The original paper measured the gender bias  
294 in Spanish and French, and proposed a mitigation method. We measure the gender bias for Romanian  
295 embeddings and provide this measure as a baseline to be improved by contributors. More specifically,  
296 we employ two sets  $(A, B)$  containing paired words that define a semantic gender direction like  
297  $(tată, mamă)$ ,  $(fiu, fiică)$ <sup>5</sup>. We also employ two sets of (unpaired) frequently-used feminine and  
298 masculine nouns to define a grammatical gender direction. The correlation between these directions  
299 is higher in Romanian (0.53) compared to the reported value in [45] for Spanish (0.39). We project  
300 the grammatical gender component into the semantic gender direction to obtain orthogonal directions.  
301 To measure the gender bias, we consider two sets  $(X, Y)$  of paired occupational embeddings, e.g.  
302  $(profesor, profesoară)$ ,  $(inginer, ingineră)$ <sup>6</sup>, etc.; see Figure 2. Using the modified WEAT metric as  
303 in [45], we compute  $b_w = ||s(w_m, A, B) - s(w_f, A, B)||$ , where  $(w_m, w_f)$  are pairs in  $(X, Y)$ , and  
304 summing  $b_w$  over the entire sets we get 2.57 (higher means more biased embeddings). This value is  
305 in between the values reported by [45] for Spanish and French (3.69 and 2.34, respectively). Our  
306 repository contains all the lists of words and a notebook to replicate this measurement.

#### 307 5.6 Error analysis

308 The newly-created datasets allow analysing the errors made by the deep models, providing a useful  
309 glimpse into how much these models capture the semantics of the Romanian language.

310 **Semantic textual similarity.** We investigate the sentence pairs from the RO-STS test set where the  
311 models make large errors, i.e. they grossly overestimate or underestimate the similarity. We consider  
312 that the error is significantly large if the difference between the ground truth and the predicted

<sup>5</sup>English *(father, mother)*, *(son, daughter)*

<sup>6</sup>English *(professor, engineer)*



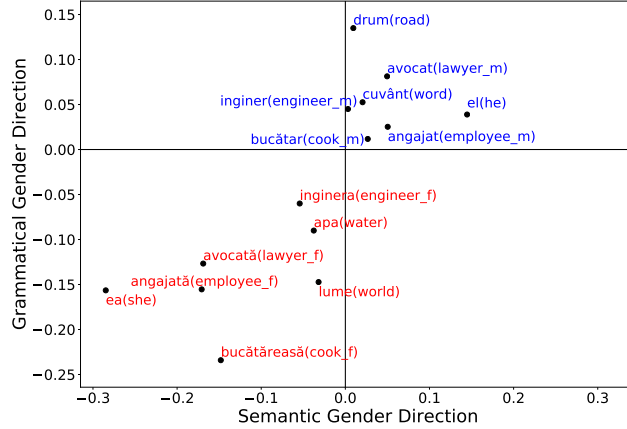


Figure 2: Occupational pairs and inanimate nouns projected on semantic gender direction (x-axis) and grammatical gender direction (y-axis). It can be observed that some feminine occupational words are farther away from the feminine definitional words than the masculine are from the masculine definitional words, revealing gender bias encoded in the embeddings.

313 similarity score is larger than an absolute value of 0.3. In this large error regime, we observe that  
 314 both models have a tendency to overestimate the similarity of sentences: 10.7% pairs for ro-BERT  
 315 and 10.6% pairs for mBERT. Moreover, mBERT has a slightly higher tendency to underestimate the  
 316 similarity compared to ro-BERT: 3.2% pairs for mBERT compared to 2.6% for ro-BERT. At closer  
 317 inspection, we observe that in most of the cases where the models overestimate the similarity, the  
 318 sentence pairs have some parts in common, either the subject, the action, or the action’s object. In  
 319 this case, the models behave similarly to a bag-of-words model. The cases where the similarity is  
 320 grossly underestimated contain idioms or the sentences have different word order. We include a few  
 321 representative examples in Table 6.

322 **Machine translation.** We manually inspect the test samples with large errors and observe that in  
 323 many cases the predicted translations are not identical, but are semantically similar to the ground  
 324 truth; for example words replaced with their synonyms (e.g. Romanian: *acum* → *în prezent*, En:  
 325 *now* → *in this moment*), adding or removing the article of a noun (e.g. Ro: *el cântă la chitară* →  
 326 *el cântă la o chitară*, En: *He is playing guitar* → *He is playing the guitar*), or even paraphrasing  
 327 entire chunks (e.g. Romanian: *Într-un sondaj realizat săptămâna trecută de CNN/ORC* → *Într-un*  
 328 *sondaj CNN/ORC de săptămâna trecută*, English: *In a poll conducted last week by CNN/ORC* →  
 329 *In a CNN/ORC poll of last week*). Such mistakes should not be penalized by a performance metric,  
 330 as it is the case with BLEU. We believe that a dataset like RO-STS might prove useful for better  
 331 estimating the quality of Romanian translations.

## 332 6 Conclusions

333 We proposed *LiRo*, a platform for benchmarking machine learning models across ten language  
 334 understanding tasks for Romanian, with the explicit goal to increase accessibility and standardization,  
 335 and to eventually accelerate progress. Additionally, we introduce, as part of *LiRo*, three new  
 336 datasets: RO-STS, XQuAD-ro, and Wiki-ro. Wiki-ro is meant to provide a standardized evaluation  
 337 dataset for language modelling. RO-STS and XQuAD-ro were obtained by human-translating their  
 338 English counterparts and represent the first datasets of their kind for the Romanian language. We  
 339 believe they play a dual role: first as standard benchmarks for Romanian semantic similarity and  
 340 Q&A, respectively, allowing the evaluation of systems dedicated to these tasks. Second, as part of  
 341 parallel corpora, they enable multilingual and cross-lingual research, which is of interest for the  
 342 wider NLP community. *LiRo* also includes tasks on cross-domain splits of the standard UD-RRR  
 343 dataset to test robustness of existing models and a task related to gender debiasing of Romanian  
 344 language embeddings, to acknowledge the importance of this line of research and encourage works  
 345 on Romanian embeddings debiasing. We pledge to continue extending *LiRo* by adding more tasks  
 346 and datasets, either by creating them from scratch or, when possible, by translating existing datasets  
 347 additionally producing parallel corpora.

348 **References**

- 349 [1] Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. On the Cross-lingual Transferability  
350 of Monolingual Representations. In *Proceedings of the 58th Annual Meeting of the Association for*  
351 *Computational Linguistics*, pages 4623–4637. Association for Computational Linguistics.
- 352 [2] Verginica Barbu Mititelu, Radu Ion, Radu Simionescu, Elena Irimia, and Ceneal-Augusto Perez.  
353 2016. The Romanian treebank annotated according to universal dependencies. In *Proceedings of the*  
354 *Tenth International Conference on Natural Language Processing*.
- 355 [3] Ondřej Bojar, Yvette Graham, Amir Kamran, and Miloš Stanojević. 2016. Results of the wmt16  
356 metrics shared task. In *Proceedings of the First Conference on Machine Translation: Volume 2,*  
357 *Shared Task Papers*, pages 199–231.
- 358 [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
359 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-  
360 Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey  
361 Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray,  
362 Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever,  
363 and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information*  
364 *Processing Systems*.
- 365 [5] Andrei Butnaru and Radu Tudor Ionescu. 2019. MOROCO: The Moldavian and Romanian  
366 dialectal corpus. In *Proceedings of the 57th Annual Meeting of the Association for Computational*  
367 *Linguistics*, pages 688–698, Florence, Italy. Association for Computational Linguistics.
- 368 [6] Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. SemEval-  
369 2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In  
370 *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages  
371 1–14, Vancouver, Canada. Association for Computational Linguistics.
- 372 [7] Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and  
373 Tony Robinson. 2013. One billion word benchmark for measuring progress in statistical language  
374 modeling.
- 375 [8] Alina Maria Ciobanu and Liviu P. Dinu. 2016. A computational perspective on the Romanian  
376 dialects. In *Proceedings of the Tenth International Conference on Language Resources and Evalua-*  
377 *tion (LREC’16)*, pages 3281–3285, Portorož, Slovenia. European Language Resources Association  
378 (ELRA).
- 379 [9] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek,  
380 Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020.  
381 Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual*  
382 *Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for  
383 Computational Linguistics.
- 384 [10] Alexis Conneau and Douwe Kiela. 2018. SentEval: An evaluation toolkit for universal sentence  
385 representations. In *Proceedings of the Eleventh International Conference on Language Resources*  
386 *and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- 387 [11] Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou.  
388 2017. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.
- 389 [12] Eugen Coşeriu. 1988. *Der romanische Sprachtypus. Versuch einer neuen Typologisierung der*  
390 *romanischen Sprachen.*, chapter IV. Gunter Narr Verlag.
- 391 [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training  
392 of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the North*  
393 *American Chapter of the Association for Computational Linguistics (NAACL)*, pages 4171–4186.

- 394 [14] Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard  
395 Socher, and Byron C Wallace. 2020. ERASER: A Benchmark to Evaluate Rationalized NLP Models.  
396 In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages  
397 4443–4458.
- 398 [15] Stefan Dumitrescu, Andrei-Marius Avram, and Sampo Pyysalo. 2020. The birth of Romanian  
399 BERT. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4324–  
400 4328, Online. Association for Computational Linguistics.
- 401 [16] Stefan Daniel Dumitrescu and Andrei-Marius Avram. 2019. Introducing RONEC—the Romanian  
402 Named Entity Corpus. *arXiv preprint arXiv:1909.01247*.
- 403 [17] Mihaela Găman and Radu Tudor Ionescu. 2020. The unreasonable effectiveness of machine  
404 learning in moldavian versus romanian dialect identification. *arXiv preprint arXiv:2007.15700*.
- 405 [18] Marcos Garcia, Carlos Gomez-Rodriguez, and Miguel A. Alonso. 2018. New treebank or repur-  
406 posed? on the feasibility of cross-lingual parsing of romance languages with universal dependencies.  
407 *Natural Language Engineering*, 24(1):91–122.
- 408 [19] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M.  
409 Wallach, Hal Daumé III, and Kate Crawford. 2018. Datasheets for datasets. *arXiv preprint*  
410 *arXiv:1803.09010*.
- 411 [20] Daniela Gerz, Ivan Vulić, Edoardo Maria Ponti, Roi Reichart, and Anna Korhonen. 2018. On  
412 the relation between linguistic typology and (limitations of) multilingual language modeling. In  
413 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages  
414 316–327, Brussels, Belgium. Association for Computational Linguistics.
- 415 [21] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. DeBERTa: Decoding-  
416 enhanced BERT with Disentangled Attention. In *International Conference on Learning Representa-*  
417 *tions*.
- 418 [22] Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson.  
419 2020. XTREME: A Massively Multilingual Multi-task Benchmark for Evaluating Cross-lingual  
420 Generalisation. In *International Conference on Machine Learning*, pages 4411–4421. PMLR.
- 421 [23] Johannes Kabatek and C. Pusch. 2015. The romance languages. In *The Languages and Linguis-*  
422 *tics of Europe*. De Gruyter.
- 423 [24] Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M Rush. 2017.  
424 OpenNMT: Open-Source Toolkit for Neural Machine Translation. In *Proceedings of the 55th Annual*  
425 *Meeting of the Association for Computational Linguistics (System Demonstrations)*, pages 67–72.
- 426 [25] Taku Kudo. 2018. Subword regularization: Improving neural network translation models with  
427 multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for*  
428 *Computational Linguistics (Volume 1: Long Papers)*, pages 66–75.
- 429 [26] Jouko Lindstedt. 2014. Balkan slavic and balkan romance: From congruence to convergence. In  
430 Juliane Besters-Dilger, Cynthia Dermarkar, Stefan Pfänder, and Achim Rabus, editors, *Congruence*  
431 *in Contact-Induced Language Change*, *Linguae & Litterae*, pages 168–183. de Gruyter, Germany.
- 432 [27] Xiaodong Liu, Kevin Duh, Liyuan Liu, and Jianfeng Gao. 2020. Very Deep Transformers for  
433 Neural Machine Translation. *arXiv e-prints*, page arXiv:2008.07772.
- 434 [28] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike  
435 Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine transla-  
436 tion. *arXiv preprint arXiv:2001.08210*.
- 437 [29] Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. The Natural  
438 Language Decathlon: Multitask Learning as Question Answering. *arXiv preprint arXiv:1806.08730*.

- 439 [30] Sungjoon Park, Jihyung Moon, Sungdong Kim, Won Ik Cho, Jiyeon Han, Jangwon Park,  
440 Chisung Song, Junseong Kim, Yongsook Song, Taehwan Oh, Joohong Lee, Juhyun Oh, Sungwon  
441 Lyu, Younghoon Jeong, Inkwon Lee, Sangwoo Seo, Dongjun Lee, Hyunwoo Kim, Myeonghwa  
442 Lee, Seongbo Jang, Seungwon Do, Sunkyoung Kim, Kyungtae Lim, Jongwon Lee, Kyumin Park,  
443 Jamin Shin, Seonghyun Kim, Lucy Park, Alice Oh, Jung-Woo Ha, and Kyunghyun Cho. 2021. Klue:  
444 Korean language understanding evaluation. *arXiv preprint arXiv:2105.09680*.
- 445 [31] Adam Poliak, Aparajita Haldar, Rachel Rudinger, J Edward Hu, Ellie Pavlick, Aaron Steven  
446 White, and Benjamin Van Durme. 2018. Collecting Diverse Natural Language Inference Problems for  
447 Sentence Representation Evaluation. In *Proceedings of the 2018 Conference on Empirical Methods*  
448 *in Natural Language Processing*, pages 67–81.
- 449 [32] Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza:  
450 A python natural language processing toolkit for many human languages. *CoRR*, abs/2003.07082.
- 451 [33] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+  
452 questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empir-*  
453 *ical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for  
454 Computational Linguistics.
- 455 [34] Georg Rehm, Maria Berger, Ela Elsholz, Stefanie Hegele, Florian Kintzel, Katrin Marheinecke,  
456 Stelios Piperidis, Miltos Deligiannis, Dimitris Galanis, Katerina Gkirtzou, Penny Labropoulou,  
457 Kalina Bontcheva, David Jones, Ian Roberts, Jan Hajič, Jana Hamrlová, Lukáš Kačena, Khalid  
458 Choukri, Victoria Arranz, Andrejs Vasiljevs, Orians Anvari, Andis Lagzdīns, Jūlija Melņika, Gerhard  
459 Backfried, Erinc Dikici, Miroslav Janosik, Katja Prinz, Christoph Prinz, Severin Stampfer, Dorothea  
460 Thomas-Aniola, José Manuel Gómez-Pérez, Andres Garcia Silva, Christian Berrío, Ulrich Germann,  
461 Steve Renals, and Ondrej Klejch. 2020. European Language Grid: An Overview. In *Proceedings*  
462 *of the 12th Language Resources and Evaluation Conference*, pages 3366–3380, Marseille, France.  
463 European Language Resources Association.
- 464 [35] Piotr Rybak, Robert Mroczkowski, Janusz Tracz, and Ireneusz Gawlik. 2020. KLEJ: Comprehen-  
465 sive Benchmark for Polish Language Understanding. In *Proceedings of the 58th Annual Meeting of*  
466 *the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 1191–1201.  
467 Association for Computational Linguistics.
- 468 [36] Julian Salazar, Davis Liang, Toan Q. Nguyen, and Katrin Kirchhoff. 2020. Masked language  
469 model scoring. In *Proceedings of the 58th Annual Meeting of the Association for Computational*  
470 *Linguistics*, pages 2699–2712, Online. Association for Computational Linguistics.
- 471 [37] Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. 2020. Green AI. *Communication*  
472 *of the ACM*, 63(12):54–63.
- 473 [38] Anca Maria Tache, Mihaela Gaman, and Radu Tudor Ionescu. 2021. Clustering Word Embed-  
474 dings with Self-Organizing Maps. Application on LaRoSeDa – A Large Romanian Sentiment Data  
475 Set. In *Proceedings of the European Chapter of the Association for Computational Linguistics*.
- 476 [39] O.M. Tomic. 2006. *Balkan Sprachbund Morpho-Syntactic Features*. Studies in Natural Language  
477 and Linguistic Theory. Springer Netherlands.
- 478 [40] Diana Trandabat, Elena Irimia, Verginica Barbu Mititelu, Dan Cristea, and Dan Tufiş. 2012. The  
479 Romanian Language in the Digital Era. *Springer, Metanet White Paper Series*.
- 480 [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
481 Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural*  
482 *Information Processing Systems*.
- 483 [42] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill,  
484 Omer Levy, and Samuel R Bowman. 2019. SuperGLUE: A stickier benchmark for general-purpose  
485 language understanding systems. In *Advances in Neural Information Processing Systems*, volume 32.

- 486 [43] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman.  
 487 2018. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding.  
 488 In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural*  
 489 *Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- 490 [44] Bryan Wilie, Karissa Vincentio, Genta Indra Winata, Samuel Cahyawijaya, Xiaohong Li,  
 491 Zhi Yuan Lim, Sidik Soleman, Rahmad Mahendra, Pascale Fung, Syafri Bahar, and Ayu Purwarianti.  
 492 2020. IndoNLU: Benchmark and resources for evaluating Indonesian natural language understanding.  
 493 In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational*  
 494 *Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages  
 495 843–857, Suzhou, China. Association for Computational Linguistics.
- 496 [45] Pei Zhou, Weijia Shi, Jieyu Zhao, Kuan-Hao Huang, Muhao Chen, Ryan Cotterell, and Kai-Wei  
 497 Chang. 2019. Examining gender bias in languages with grammatical gender. In *Proceedings of the*  
 498 *2019 Conference on Empirical Methods in Natural Language Processing and the 9th International*  
 499 *Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5276–5284, Hong  
 500 Kong, China. Association for Computational Linguistics.

## 501 Checklist

- 502 1. For all authors...
- 503 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 504 contributions and scope? [Yes]
- 505 (b) Did you describe the limitations of your work? [Yes] The newly created datasets are  
 506 suitable mainly for evaluation, as their size is relatively small; the details are included  
 507 in Section 4. We hope to propose larger datasets in the future, suitable also for training.
- 508 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Our work  
 509 is aimed at accelerating research and implicitly the adoption of NLP solutions for  
 510 Romanian language. Such tools come with potential issues (e.g. biases) which can  
 511 have negative impact. We briefly discuss this in section 5.5 and propose as part of the  
 512 benchmark a task to measure and mitigate gender bias in learnt embeddings.
- 513 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 514 them? [Yes]
- 515 2. If you are including theoretical results...
- 516 (a) Did you state the full set of assumptions of all theoretical results? [N/A]  
 517 (b) Did you include complete proofs of all theoretical results? [N/A]
- 518 3. If you ran experiments (e.g. for benchmarks)...
- 519 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
 520 mental results (either in the supplemental material or as a URL)? [Yes] As URLs.
- 521 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 522 were chosen)? [Yes]
- 523 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
 524 ments multiple times)? [N/A]
- 525 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 526 of GPUs, internal cluster, or cloud provider)? [N/A]
- 527 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 528 (a) If your work uses existing assets, did you cite the creators? [Yes]  
 529 (b) Did you mention the license of the assets? [Yes]  
 530 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
 531 As URLs.
- 532 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
 533 using/curating? [N/A]  
 534 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
 535 information or offensive content? [N/A] It does not contain personally identifiable  
 536 information or offensive content.

- 537 5. If you used crowdsourcing or conducted research with human subjects...
- 538 (a) Did you include the full text of instructions given to participants and screenshots, if
- 539 applicable? [Yes] See Section 4, translation of RO-STS using volunteers.
- 540 (b) Did you describe any potential participant risks, with links to Institutional Review
- 541 Board (IRB) approvals, if applicable? [N/A]
- 542 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 543 spent on participant compensation? [No] The annotators on RO-STS were volunteers.
- 544 The translation of XQUAD-ro was done through a professional translation operator,
- 545 not a crowdsourcing platform.