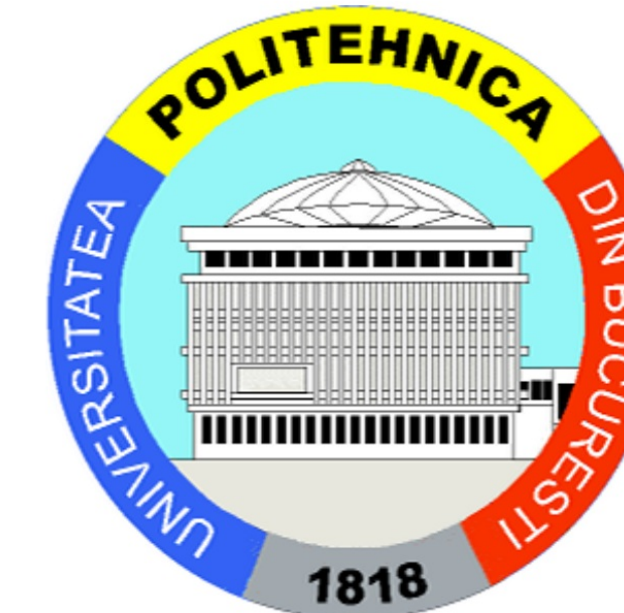


DISTILLING THE KNOWLEDGE OF ROMANIAN BERTS USING MULTIPLE TEACHERS

Andrei-Marius Avram¹, Darius Catrina¹, Dumitru-Clementin Cercel², Mihai Dascalu², Traian Rebedea², Vasile Păiș¹, Dan Tufiș¹

¹Research Institute for Artificial Intelligence, Romanian Academy,

²University Politehnica of Bucharest, Faculty of Automatic Control and Computers



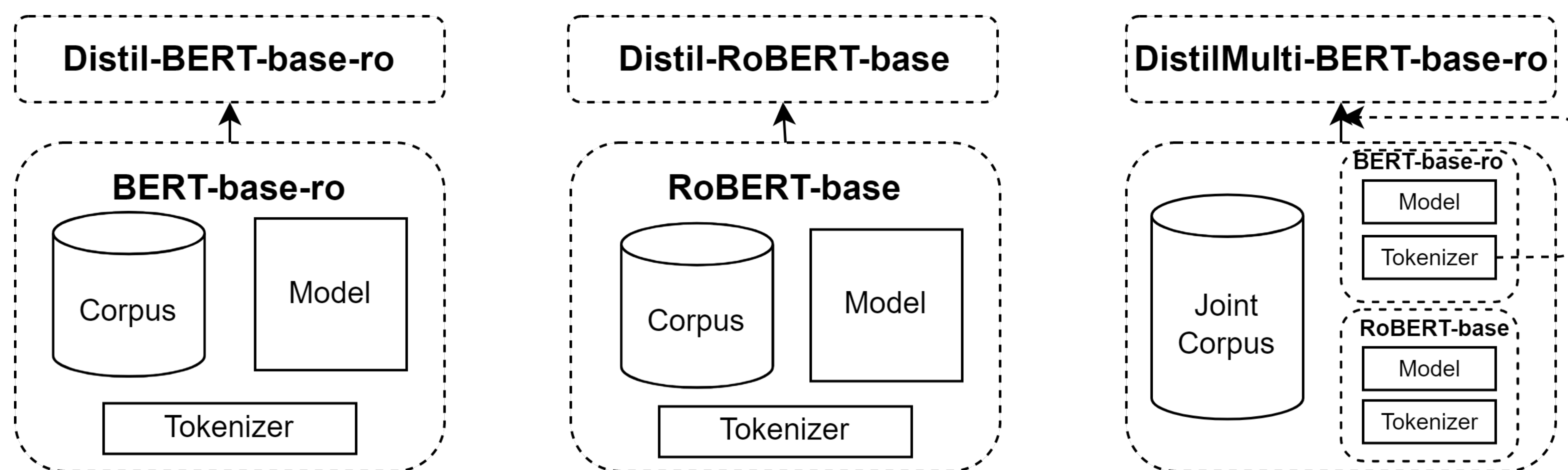
Motivation

- Running large-scale pre-trained language models in computationally constrained environments is a challenging problem
- Model optimizations (e.g., distillation, quantization) focus mostly on the English language, thus widening the gap when considering low-resource languages
- As a response, we introduce and evaluate three distilled models for the Romanian language:
 1. Distil-BERT-base-ro
 2. Distil-RoBERT-base
 3. DistilMulti-BERT-base-ro

Task Evaluation

- *Part-of-Speech Tagging*: Universal Part-of-Speech (UPOS) and eXtended Part-of-Speech (XPOS)
- *Named Entity Recognition (NER)*: Tag a sequence of tokens with IOB entities
- *Sentiment Analysis*: Positive or negative sentiment (SAPN) or rating (SAR)
- *Dialect Identification (DI)*: Identify the Romanian/Moldavian dialects in news articles
- *Semantic Textual Similarity (STS)*: Given a pair of sentences, measure their semantically relatedness

Distillation Process



Results

Model	UPOS	XPOS	NER	SAPN	SAR	DI	STS
BERT-base-ro	98.00	96.46	85.88	97.94	79.61	95.41	79.94
RoBERT-base	98.02	97.15	85.14	98.20	79.40	96.17	81.18
<i>Distil-BERT-base-ro</i>	97.97	97.08	83.35	98.12	80.51	96.31	80.57
<i>Distil-RoBERT-base</i>	97.12	95.79	83.11	97.61	79.58	96.11	79.80
<i>DistilMulti-BERT-base-ro</i>	98.07	96.93	83.22	97.74	79.77	96.18	80.66

Loyalty

Regression loyalty (R-L): the Pearson Coefficient between the outputs of the teacher and the student

Model	L-L	P-L	R-L
BERT-base-ro	87.80	74.76	94.96
RoBERT-base	84.63	71.95	92.24
Distil-BERT-base-ro	89.75	73.23	94.64

Conclusions

- Introduced three distilled models for the Romanian language that are approximately 35% smaller and that maintain most of the performance of their teachers
- Evaluated these distilled models on five Romanian tasks in comparison to their original teachers
- Introduced a new metric that traces how well the student and teacher predictions are aligned in the regression tasks – *regression loyalty*