# Benchmarking Multilabel Topic Classification in the Kyrgyz Language

Anton Alekseev[1,2,3,4][0000−0001−6456−3329],
Sergey Nikolenko[1,2,3][0000−0001−7787−2251], and
Gulnara Kabaeva[4][0000−0003−3001−7201]

[1] Steklov Mathematical Institute at St. Petersburg, St. Petersburg, Russia
[2] St. Petersburg State University, St. Petersburg, Russia
[3] Kazan (Volga Region) Federal University, Kazan, Russia
[4] Kyrgyz State Technical University n. a. I. Razzakov, Bishkek, Kyrgyzstan

**Abstract.** Kyrgyz is a very underrepresented language in terms of modern natural language processing resources. In this work, we present a new public benchmark for topic classification in Kyrgyz, introducing a dataset based on collected and annotated data from the news site *24.KG* and presenting several baseline models for news classification in the multilabel setting. We train and evaluate both classical statistical and neural models, reporting the scores, discussing the results, and proposing directions for future work.

**Keywords:** Topic classification · Kyrgyz language · Multi-label classification · Low-resource languages.

## 1 Introduction

Kyrgyz is an agglutinative Turkic language spoken in several countries, notably China and Tajikistan in addition to Kyrgyzstan; it is by no means an endangered language, and several millions of people call it their mother tongue [35]. However, despite a large amount of linguistic work, including computational linguistics (see Section 2), it is certainly a *low-resource* language, with a very modest number of tools and datasets available in the open for Kyrgyz language processing[5]. A recent publication [35], following the taxonomy proposed in [22], labels Kyrgyz with the "Scraping By" status, defined as follows: "With some amount of unlabeled data, there is a possibility that they could be in a better position in the 'race' in a matter of years. However, this task will take a solid, organized movement that increases awareness about these languages, and also sparks a strong effort to collect labelled datasets for them, seeing as they have almost none." Therefore, we believe that a meaningful effort to construct open manually annotated text collections or other reliable resources for Kyrgyz

---

[5] For a list of tools, corpora, and other language resources for Turkic languages including Kyrgyz, see e.g. https://github.com/alexeyev/awesome-kyrgyz-nlp and http://ddi.itu.edu.tr/en/toolsandresources.

language processing is in great demand; modern NLP, while shifting towards universal models, is still hard to imagine without at least evaluation data.

Text topic classification is a core task in natural language processing and information retrieval [34]; it is one of the most popular in practice, with applications in advertising [71], news aggregation, and many other industries. Very often, topic categorization is posed as a multilabel classification problem since the same text can touch upon multiple topics [33, 46, 69].

In this work, we present virtually the first labeled dataset for text classification in the Kyrgyz language based on the *24.kg* news portal. Moreover, we propose several baseline models and evaluate their results; in this evaluation, we see that multilingual models do help to process Kyrgyz, using even primitive stemming and passing from word $n$-grams to symbol $n$-grams quite expectedly help, and deep learning models that we have considered perform better than the best linear models with virtually no hyperparameter search. Thus, the contributions of our work are threefold: (1) a novel manually labeled dataset for texts in the Kyrgyz language, (2) several approaches to multi-label Kyrgyz text classification, (3) proof of concept for the feasibility of multilingual LLMs for Kyrgyz language processing in supervised tasks. The paper is organized as follows: Section 2 discusses related work, Section 3 introduces our dataset, Section 4 shows baseline models and experimental setup used in our experiments, Section 5 discusses experimental results, and Section 6 concludes the paper.

## 2 Related Work

*Topic classification.* Text topic classification is one of the oldest and best known tasks in information retrieval and natural language processing [34, 72]. It is straightforwardly defined as a supervised learning (classification) task, often expanding into multilabel classification since longer texts are hard to fit into a single topic [33, 46, 69], a problem that is still attracting attention in the latest deep learning context [32, 62]. Many approaches have been developed for datasets of different nature: (i) news article datasets such as BBC News [18], Reuter [30], 20 Newsgroups [28], or WMT News Crawl [29]; (ii) scientific texts such as arXiv abstracts [29], patents [56], or clinical texts [42, 53, 66]; (iii) social media posts where topics are usually represented by [hash]tags [13], and more. We note especially prior efforts related to text classification for low-resource languages [2, 11, 14, 16, 19].

*Kyrgyz language processing.* There already exists a large corpus of linguistic research papers dedicated to various aspects of the Kyrgyz language: (1) grammar, syntax and morphology modeling [7, 21, 24–26, 50, 57, 64, 65, 68, 74–77, 79, 82], including a recent release of 780 dependency trees [5] as part of the Universal Dependencies initiative [41], (2) text-related statistics and construction of corpora/dictionaries [4, 24–26, 73, 80, 83], (3) computer-aided language learning and other educational systems [9, 23], (4) machine translation [45, 58, 64, 78], (5) lexicons and thesauri [7, 84], (6) computational linguistics in general [81], and more.

Kyrgyz also appears in multiple works as a part of multilingual research studies, e.g. on multiway machine translation [36, 37] and even text categorization [31], although the latter uses a different (Arabic) script than our work.

News articles represent a traditional and widely used text domain, traditionally a valuable source of data both for information retrieval and natural language processing. News-based datasets have found many applications including such non-traditional ones as relation classification [49]. In this work, we concentrate on the news domain primarily because to the best of our knowledge, for the Kyrgyz language only fiction, news, and Wikipedia articles are readily available online. Collecting social media content, which may also be useful for numerous NLP tasks [3, 8, 27, 38–40, 59–61], is also possible but requires significant extra effort for preprocessing; in particular, preprocessing steps for social media sources would have to include language detection since oftentimes people writing in Kyrgyz also publish posts in Russian and other languages.

Several research efforts on Kyrgyz open corpora and dictionaries are currently in progress (see, e.g., [24–26]), but as of 2023, there are still very few manually annotated datasets useful for Kyrgyz language processing. We hope to start filling this gap with this work.

## 3   Dataset

### 3.1   Annotation

With permission of *24.kg*[6] editors, we collected 23 283 news articles in Kyrgyz, dated from May 2017 to October 2022. The portal does not provide any topical tags for articles in Kyrgyz, hence we had to either match collected articles with possibly available articles in Russian, which are tagged, or annotate them with our own topical categories. The original rubrics used at *24.kg* include: (1) Власть (government, politics and law), (2) Общество (society), (3) Экономика (economics), (4) Происшествия (accidents, current events), (5) Агент 024 (current events), (6) Спорт (sports), (7) Техноблог (tech), (8) Спецпроекты (special projects), (9) Кыргызча (articles in Kyrgyz), (10) English (articles in English), (11) Бизнес (business). Some of the rubrics are clearly not topical ("English", "Кыргызча"), some are multi-topic ("Спецпроекты", "Агент 024"), and some other topics also turned out to cover very diverse information. Therefore, we had to introduce our own topical labels.

While some general-purpose taxonomies for content classification do exist and are used, e.g., in advertising, including *dmoz*[7] and IAB[8] taxonomies, the label sets there are too broad for the purpose of news classification. Our preliminary experiments with translated news titles zero-shot classification with IAB Tier 1 tags (in the *label-fully-unseen* setting [67]) yielded poor prediction quality.

---

[6] https://24.kg/

[7] https://www.dmoz-odp.org/; previously https://www.dmoz.org/.

[8] https://iabtechlab.com/standards/content-taxonomy/

| Title | Proposed labels |
|-------|-----------------|
| The presidential candidate who violated traffic rules paid... | law/crime, politics |
| Cars of drivers who do not pay fines on time will be... | law/crime |
| 44 percent of the 108,000 fines imposed for violating traffic... | law/crime |
| Party candidate was fined 7,500 soms for holding a concert... | law/crime, politics |
| Fines for garbage thrown from cars have been increased... | law/crime, ecology |

**Table 1.** Sample cluster (#16).

However, we still consider this direction very promising from the practical point of view and leave it for future research.

To motivate the introduction of a custom set of labels, we have automatically translated[9] titles of Kyrgyz articles into English, randomly sampled 500 out of 23 284 of them (a subset small enough to annotate in reasonable time yet hopefully large enough to derive meaningful conclusions regarding the topics), and obtained their embeddings via the SentenceBERT model [48] (`all-mpnet-base-v2`, the best-performing[10] fine-tuned MPNet model [54]). Then, we have grouped the resulting embeddings using agglomerative clustering (Euclidean distance, Ward linkage [63], other hyperparameters left at default as provided by *scikit-learn* version 1.0 [44]) into 100 clusters. Note that the exact clustering procedure and chosen parameters are not of significant importance here; the main idea is to group texts into hopefully small clusters of very similar titles to speed up annotation and, most importantly, to be able to easily invent topic names that are neither too general nor too specific. Note also that we had to translate the titles and apply the model trained on English language data not for the annotation itself but only because there are no good sentence embeddings models for Kyrgyz with reported quality. Where it was impossible to deduce the topics of the article from the title, we made decisions based on the original Kyrgyz news texts. A sample cluster is presented in Table 1.

The exploratory annotation task was defined as follows: for each cluster, invent a topic name that best describes most if not all titles and use it as the class label. Then correct the label for titles in the cluster that do not fit the invented topic. If multiple topics apply to some of the titles, add more tags where necessary. After that, we make another pass over all 500 titles since some of the labels were not "available" at the beginning, i.e., a certain "general" label might be added to the label set after some of the texts that would be appropriately labeled with it had already been annotated. As a result, we obtained a refined list of 20 labels shown in Table 2. To validate the label set by comparing label distributions with each other, we have annotated 500 more English translations of the titles using the same set of labels. We found that the difference in label count distributions in the two sets of 500 were relatively small, which showed that our annotation was consistent. Finally, we have annotated 500 more texts

---

[9] Google Translate: https://translate.google.com/?sl=ky&tl=en&op=docs

[10] As of 24.08.2023: https://www.sbert.net/docs/pretrained_models.html.

| Class label | 500-1 | 500-2 | Description |
|---|---|---|---|
| politics | 127 | 174 | Mentions of politicians and political decisions |
| law/crime | 126 | 128 | Judiciary and penitentiary sys., legislature, trials, crime |
| foreign affairs | 84 | 91 | Any non-Kyrgyzstan-related news |
| health | 68 | 63 | Health and medicine-related news (mostly COVID19) |
| local | 43 | 43 | Traffic rerouting, events scheduling |
| accidents | 41 | 31 | Disasters, fires, road accidents, etc. |
| econ/finance | 37 | 50 | Money, import-export and labour-related news |
| society | 36 | 49 | Local initiatives, protests, other citizen-related news |
| culture | 32 | 29 | Cultural events and initiatives, celebrity news |
| citizens abroad | 17 | 11 | Migration questions and Kyrgyz people abroad |
| sports | 16 | 15 | Awards, announcements, famous sportspeople mentions |
| natural hazards | 13 | 5 | Inconveniences and threats due to natural reasons |
| development | 12 | 25 | Realty, land use and infrastructural development |
| religion | 11 | 13 | Religion-related news |
| science/tech | 9 | 7 | Everything related to science and technology |
| border | 8 | 9 | Kyrgyzstan's borders-related conflicts and resolutions |
| education | 7 | 22 | News on educational procedures/events/institutions |
| weather | 6 | 4 | Weather forecasts and reports |
| ecology | 4 | 4 | Ecological initiatives, laws and reports |
| natural resources | 2 | 0 | Issues related to natural resources |

**Table 2.** Topical tags for *24.kg*, total number of tags in the first two annotation batches and their descriptions. The counts do not sum up to 500 since this is a multilabel task.

| | # texts | # sent. | Sent/text | # tokens | Unique tokens | Tok/text | Unique stems |
|---|---|---|---|---|---|---|---|
| **Train** | 1 000 | 7 319 | 7.32 ± 5.36 | 107 556 | 18 958 | 107.56 ± 74.02 | 9 872 |
| **Test** | 500 | 4 025 | 8.05 ± 8.78 | 57 414 | 12 885 | 114.83 ± 101.15 | 6 924 |

**Table 3.** Dataset statistics; sentences counted via the `sent_tokenize` method from NLTK [6]. Per-text values show mean value and standard deviation.

following the same procedure. We do not disclose the exact distribution of the labels in the final batch for the sake of fairness in possible future competitions: knowing the exact number of texts with a certain label might be used as a test data leak to improve results.

## 3.2 Data Description

The dataset consists of 1 500 texts, annotated in three sessions as described above. Since the dataset is relatively small, we split it in only two parts: the first two batches of 500 (i.e., *training* set has 1 000 texts) and the last batch (i.e. the *test* set has 500 texts).

For further application of models based on the bag-of-ngrams approach, the texts had to be split into tokens and, possibly, stemmed/lemmatized. For tokenization, we used the splitting mechanism provided by the Apertium Project morphological analyzer [17, 65]; to the best of our knowledge, this is the only open source engine for Kyrgyz morphology. Similarly, for *word normalization* we used the Apertium-Kir [65] FST's token segmentation; since prefixes are uncommon in Kyrgyz, the first segment was used as the stem. Overall dataset statistics are presented in Table 3.

## 4   Models and Experimental Setup

The resulting dataset is far too small to be used for training, especially for classical models that heavily depend on frequency estimates of various ratios of tokens and n-grams, e.g., models based on the bag-of-words assumption. However, we can use the dataset in cross-validation to make comparisons across models that perform transfer learning. Still, we include classical approaches into the benchmark as well, since several works have demonstrated that word/character n-gram baselines are sometimes surprisingly competitive, e.g., in entity linking [1,51], so they should not be ignored even for a relatively small dataset.

We have used grid search to find the best parameters. Since the training set is small and imbalanced in terms of labels, we used 2-fold validation for hyperparameter search with a stratified split into two subsets preserving the label distribution[11]. Below we show the considered values and ranges of hyperparameters in addition to the models themselves.

### 4.1   Approaches based on the bag-of-ngrams assumption

To provide a classic baseline, we have considered several sparse text representations (essentially bags-of-ngrams) and several corresponding models.

*Text preprocessing.* We tested several text representations. First, we tokenized text (Section 3.2) into unigrams, 1-2-grams, 1-2-3-grams, and 2-3-grams. For frequency cutoffs we retained tokens with maximum document frequency (maxdf) of 40%, 60%, 80%, and 100% and minimum occurrence (mincount) in 2, 5, or 10 documents. We also set the maximum number of features (maxfeat) equal to 2 000 or 10 000. In another set of experiments, we used character n-grams: 2-3-grams, 3-4-grams, and 5-6-grams; maxdf for a character n-gram was set to 40%, 70%, or 100%, mincount was 4, 10, and 15, and maxfeat was 2 000 or 10 000. Then, having stemmed the texts as in Section 3.2, we have run experiments with the same "vectorization" parameters.

---

[11] Specifically, we used the *IterativeStratification* algorithm from the *scikit-multilearn* library [55].

*Independent classifiers ("Independent").* In this set of baselines, we train a separate model for every label, using models that are known to perform well for sparse features: logistic regression (both LBFGS and SGD optimization methods), a linear model with hinge loss (linear SVM), and a linear model with Huber loss (usually preferred for regression tasks). In the search for the best model, we treated log-loss, hinge loss, and Huber loss as hyperparameters. Experiments with LBFGS were carried out separately, which is reflected in the results table in Section 5. For logistic regression with the LBFGS optimizer (that includes $L_2$-regularization), we have tested the performance with regularization strength $C = \frac{1}{\lambda} \in \{0.7, 0.9, 1.0\}$ and limited the number of iterations to $1\,000$ or $10\,000$ steps. For other models, apart from the loss function, we have tested the averaging mechanism (enabling/disabling it), $L_1$, and $L_2$-regularizers, and limited the number of iterations to either $20\,000$ or $100\,000$ steps.

*Binary classifiers chain ("Chain").* In this approach, the base models (which are the same as in the previous paragraph) make predictions in a sequence; the training task for every label in a chain includes predictions for previous labels as features. Apart from the hyperparameters listed in the previous paragraph, we have tried different orderings of the prediction chain.

*Multilabel k-Nearest Neighbors.* We have added two models based on k-nearest neighbors to the grid search: (1) the model *ML-kNN* introduced in [70], which uses Bayesian inference to assign labels to test classes based on the standard kNN output, (2) a binary relevance kNN classifier (*BR-kNN*), a similar method introduced in [15] that assigns the labels that have been assigned to at least half of the neighbors. Although nearest neighbors classifiers are known to perform poorly for high-dimensional vectors (which bags-of-ngrams are), we added them to the task due to them being multi-label *by design*. We have tested $k \in \{1, 2, 3, 5, 10\}$ neighbors; for *ML-kNN*, we tried different values of the smoothing parameter, $s \in \{0.1, 0.5, 0.7, 1.0\}$. As a reliable implementation, we used a combination of models from the *scikit-learn* and *scikit-multilearn Python* libraries [44, 55].

## 4.2   Neural baseline

As a modern approach to fine-tuning neural networks, we have intentionally selected the most standard method, which is not necessarily state of the art for other languages. Among multilingual pretrained large language models, one of the most popular ones is XLM-RoBERTa (large)[12], which is essentially a RoBERTa model [10] trained on a 2.5TB segment of *CommonCrawl* data containing 100 languages, including Kyrgyz. We used XLM-RoBERTa in the multilabel classification fine-tuning setting, using a "classification head" with two linear feedforward layers with dropout and a binary cross-entropy loss. We have used the same split as before as the train-development split to find the best number of epochs (14 out of 15) based on the Jaccard score metric (see below). We used

---

[12] Available on HuggingFace: https://huggingface.co/xlm-roberta-large.

the AdamW optimizer [20] (the AMSGrad version [47]), with weight decay set
to 0.01, learning rate set to 0.00002, and exponential learning rate scheduling
with $\gamma$ coefficient set to 1.0; $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$. Batch size was
set to 4 mostly due to the equipment-related constraints. Also note that the
pretrained `xlm-roberta-large` checkpoint uses byte-pair encoding (bpe) [52] as
the tokenizer (provided with the model).

### 4.3  Evaluation metrics

Each prediction is a set of labels represented as a vector of 0s and 1s, where
1 means that the corresponding label has been predicted. Several metrics from
"regular" binary and multi-class classification can also be applied for multilabel
classification. The counterpart of accuracy here is the fraction of exact matches
("Exact"). The $F_1$ measure (a harmonic mean of precision and recall) can be
computed for every sample if we treat each binary vector representation of the
label set as binary prediction results and then averaged ("F1-sample" in Table 4).
Besides, the $F_1$ measure can be computed for each label and, e.g., micro-averaged
("F1-micro"). We also used metrics unique to the multilabel setting: share of
samples where at least one label is predicted correctly ("@l1") and the Hamming
loss ("Hamm"), i.e., the Hamming distance between binary vectors of labels.
Finally, we report the metric we have used for model selection: the sample-
averaged Jaccard similarity computed for each pair of predicted and ground
truth label sets; Jaccard similarity between sets $A$ and $B$ is defined as $\frac{|A \cap B|}{|A \cup B|}$.

## 5   Results

Results of our computational experiments are presented in Table 4. It clearly
demonstrates that employing the multilingual models for supervised tasks with
Kyrgyz text data is feasible, since a fine-tuned *XLM-RoBERTa*-based classifier
(without any hyperparameter search) outperforms all other approaches. Note
that this result has been far from obvious, since, in our preliminary experi-
ments, fine-tuning another popular model `bert-base-multilingual-cased` (in
our case essentially BERT [12] with an added feedforward layer and dropout)
did not bring any meaningful results. Another interesting observation is that
while (i) Apertium-Kir does not consider the contexts of words, (ii) it is not
a lemmatizer, and (iii) the selected stemming method is very far from perfect,
even this kind of text normalization does bring improvements compared to the
basic bag-of-ngrams approach. Moving from word ngrams to character ngrams
also improves the results in most cases, which one could expect since the Kyrgyz
language is morphologically rich.

## 6   Conclusion

In this work, we have introduced a new annotated text collection in the Kyrgyz
language for multilabel topic classification and evaluated several baseline mod-

| Configuration | JaccCV↑ | @l1↑ | Jaccard↑ | Exact↑ | Hamm↓ | F1-micro↑ | F1-sample↑ |
|---|---|---|---|---|---|---|---|
| **Bag-of-Token-Ngrams** | | | | | | | |
| Independent, LBFGS, 1-gram | .390 | .59 | .43 | .29 | .06 | .54 | .48 |
| Chain, LBFGS, 1-gram | .405 | .63 | .46 | .31 | .06 | .56 | .51 |
| Independent, SGD, hinge loss, 1-2-gram | .465 | .68 | .47 | .29 | .06 | .56 | .53 |
| Chain, SGD, hinge loss, 1-gram | .474 | .68 | .49 | .32 | .06 | .56 | .55 |
| ML-kNN, 1 neighbor, 0.1-smoothing, 1-gram | .276 | .45 | .30 | .19 | .10 | .33 | .35 |
| BRML-kNN, 1 neighbor,1-gram | .276 | .45 | .30 | .19 | .10 | .33 | .35 |
| **Bag-of-Token-Character-Ngrams** | | | | | | | |
| Independent, LBFGS, 2-3-grams | .491 | .69 | .49 | .33 | .06 | .58 | .55 |
| Chain, LBFGS, 2-3-grams | .494 | .69 | .49 | .33 | .06 | .58 | .55 |
| Independent, SGD, hinge loss, 3-4-grams | .521 | .70 | .46 | .26 | .07 | .55 | .54 |
| Chain, SGD, hinge loss, 3-4-ngrams | .524 | .71 | .48 | .28 | .07 | .55 | .55 |
| ML-kNN, 1 neighbors, 0.1-smoothing, 2-3-gram | .412 | .65 | .42 | .24 | .08 | .48 | .49 |
| BRML-kNN, 1 neighbor, 2-3-gram | .412 | .65 | .42 | .24 | .08 | .48 | .49 |
| **Bag-of-Stem-Ngrams** | | | | | | | |
| Independent, LBFGS, 1-gram | .451 | .67 | .50 | .34 | .05 | .59 | .55 |
| Chain, LBFGS, 1-gram | .463 | .68 | .51 | .35 | .05 | .60 | .56 |
| Independent, SGD, log loss, 1-gram | .514 | .74 | .52 | .33 | .06 | .61 | .59 |
| Chain, SGD, 1-gram | .516 | .74 | .54 | .36 | .06 | .61 | .60 |
| ML-kNN, 1 neighbor, 0.1-smoothing, 1-gram | .345 | .55 | .36 | .21 | .09 | .41 | .42 |
| BRML-kNN, 1 neighbor,1-gram | .345 | .55 | .36 | .21 | .09 | .41 | .42 |
| **Bag-of-Stem-Character-Ngrams** | | | | | | | |
| Independent, LBFGS, 2-4-grams | .494 | .71 | .52 | .35 | .06 | .61 | .58 |
| Chain, LBFGS, 5-6-grams | .490 | .70 | .51 | .35 | .06 | .60 | .57 |
| Independent, SGD, hinge loss, 3-4-grams | .522 | .70 | .49 | .32 | .06 | .58 | .55 |
| Chain, SGD, hinge loss, 3-4-grams | .524 | .69 | .50 | .33 | .06 | .58 | .56 |
| ML-kNN, 1 neighbor, 0.1-smoothing, 3-4-grams | .425 | .65 | .42 | .25 | .08 | .46 | .49 |
| BRML-kNN, 1 neighbor, 3-4-grams | .425 | .65 | .42 | .25 | .08 | .46 | .49 |
| XLM-RoBERTa (with bpe tokenization) | | | **.88** | **.66** | **.46** | **.04** | **.72** | **.73** |

**Table 4.** Evaluation results: @l1 — "at-least-one", Hamm — Hamming loss, JaccCV — mean Jaccard score in cross-validation, ↑ — more is better, ↓ — less is better.

els[13]. This is one of the first open datasets for the low-resource Kyrgyz language. As for baselines, we have found that while classical baselines can achieve accept-

---

able results, especially after (even primitive) stemming, a straightforward neural baseline achieves significantly better results even with virtually no hyperparameter search.

In the future, we plan to further improve the current labeling scheme, additionally expanding and validating current annotations; we plan to ask multiple experts to label the texts using models trained on currently presented data to speed up the labeling. Then, we plan to increase the dataset size by annotating more news texts. Afterwards, we plan to hold a competition that should uncover state of the art multilabel classification models for the Kyrgyz language news domain.

Also, to enhance the benchmark with an arguably even more fair comparison, we plan to: (1) translate original texts to English via *Google Translate* and report the scores of the relevant neural models that employ English LLMs as backbones or the scores of zero-shot classification via prompting state of the art generative models such as, e.g., GPT-4 [43]; (2) add the results of the *fastText* supervised classification model trained on our data to the benchmark after publication; (3) study whether using data from a similar domain in other Turkic languages can help improve classification quality. In general, we hope that the presented dataset will be able to serve as the basis for these and other experiments and become a starting point for novel NLP research for the Kyrgyz language.

# References

1. Alekseev, A., Miftahutdinov, Z., Tutubalina, E., Shelmanov, A., Ivanov, V., Kokh, V., Nesterov, A., Avetisian, M., Chertok, A., Nikolenko, S.: Medical crossing: a cross-lingual evaluation of clinical entity linking. In: Proceedings of the Thirteenth Language Resources and Evaluation Conference. pp. 4212–4220 (2022)
2. An, B.: Prompt-based for low-resource tibetan text classification. ACM Trans. Asian Low-Resour. Lang. Inf. Process. (may 2023), just Accepted
3. Apishev, M., Koltsov, S., Koltsova, O., Nikolenko, S.I., Vorontsov, K.: Mining ethnic content online with additively regularized topic models. Computación y Sistemas **20**(3), 387–403 (2016)
4. Baisa, V., Suchomel, V.: Turkic language support in sketch engine. In: Proceedings of the International conference "Turkic Languages Processing" TurkLang-2015. pp. 214–223 (2015)
5. Benli,      I.:      Ud_kyrgyz-ktmu:      Ud      for      kyrgyz      (2023), https://github.com/UniversalDependencies/UD_Kyrgyz-KTMU/
6. Bird, S., Klein, E., Loper, E.: Natural language processing with Python: analyzing text with the natural language toolkit. " O'Reilly Media, Inc." (2009)
7. Boizou, L., Mambetkazieva, D.: From kyrgyz internet texts to an xml full-form annotated lexicon: a simple semi-automatic pipeline. In: TurkLang 2017: Пятая международная конференция по компьютерной обработке тюркских языков:

Труды конференции. Т 1. Казань: Издательство Академии наук Республики Татарстан, 2017. Казань: Издательство Академии наук Республики Татарстан, 2017 (2017)

8. Buraya, K., Farseev, A., Filchenkov, A., Chua, T.S.: Towards user personality profiling from multiple social networks. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 31 (2017)

9. Cetin, M.A., Ismailova, R.: Assisting tool for essay grading for turkish language instructors. MANAS Journal of Engineering **7**(2), 141–146 (2019)

10. Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, É., Ott, M., Zettlemoyer, L., Stoyanov, V.: Unsupervised cross-lingual representation learning at scale. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 8440–8451 (2020)

11. Cruz, J.C.B., Cheng, C.: Establishing baselines for text classification in low-resource languages. CoRR **abs/2005.02068** (2020), https://arxiv.org/abs/2005.02068

12. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). pp. 4171–4186 (2019)

13. Dhingra, B., Zhou, Z., Fitzpatrick, D., Muehl, M., Cohen, W.: Tweet2Vec: Character-based distributed representations for social media. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). pp. 269–274. Association for Computational Linguistics, Berlin, Germany (Aug 2016)

14. Ein-Dor, L., Halfon, A., Gera, A., Shnarch, E., Dankin, L., Choshen, L., Danilevsky, M., Aharonov, R., Katz, Y., Slonim, N.: Active learning for BERT: An empirical study. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 7949–7962. Association for Computational Linguistics, Online (Nov 2020)

15. Eleftherios Spyromitros, Grigorios Tsoumakas, I.V.: An empirical study of lazy multilabel classification algorithms. In: Proc. 5th Hellenic Conference on Artificial Intelligence (SETN 2008) (2008)

16. Fesseha, A., Emiru, E., Diallo, M., Dahou, A.: Text classification based on convolutional neural networks and word embedding for low-resource languages: Tigrinya. Information **12**, 52 (01 2021)

17. Forcada, M.L., Ginestí-Rosell, M., Nordfalk, J., O'Regan, J., Ortiz-Rojas, S., Pérez-Ortiz, J.A., Sánchez-Martínez, F., Ramírez-Sánchez, G., Tyers, F.M.: Apertium: a free/open-source platform for rule-based machine translation. Machine translation **25**, 127–144 (2011)

18. Greene, D., Cunningham, P.: Practical solutions to the problem of diagonal dominance in kernel document clustering. In: Proceedings of the 23rd International Conference on Machine Learning. p. 377–384. ICML '06, Association for Computing Machinery, New York, NY, USA (2006)

19. Homskiy, D., Maloyan, N.: Dn at semeval-2023 task 12: Low-resource language text classification via multilingual pretrained language model fine-tuning (2023)

20. I. Loshchilov, F.H.: Decoupled weight decay regularization. In: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019 (2019)

21. Israilova, N.A., Bakasova, P.S.: Morphological analyzer of the kyrgyz language. In: Proceedings of the V International Conference on Computer Processing of Turkic Languages Turklang. vol. 2, pp. 100–116 (2017)

22. Joshi, P., Santy, S., Budhiraja, A., Bali, K., Choudhury, M.: The state and fate of linguistic diversity and inclusion in the nlp world. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 6282–6293 (2020)

23. Karabaeva, S., Dolmatova, P., Imanalieva, A.: Computer-mathematical modeling of national specificity of spatial models in kyrgyz language. In: Proceedings of the International conference "Turkic Languages Processing" TurkLang-2015. pp. 416–422 (2015)

24. Kasieva, A.A., Kadyrbekova, A.K.: Corpus annotation tools: Kyrgyz language corpus (using turkic lexicon apertium and penn treebank tools). In: Общество, язык и культура XXI века. pp. 207–214 (2021)

25. Kasieva, A.A., Satybekova, A.T.: Parts-of-speech annotation of the newly created kyrgyz corpus. Herald of KRSU **20**(6), 67–72 (2020)

26. Kasieva, A., Knappen, J., Fischer, S., Teich, E.: A new kyrgyz corpus: sampling, compilation, annotation. In: Jahrestagung der Deutschen Gesellschaft für Sprachwissenschaft (poster session) (2020), http://hdl.handle.net/21.11119/0000-0004-B62D-D

27. Koltsova, O., Koltsov, S., Nikolenko, S.I.: Communities of co-commenting in the russian livejournal and their topical coherence. Internet Research **26**(3), 710–732 (2016)

28. Lang, K.: Newsweeder: Learning to filter netnews. In: Prieditis, A., Russell, S. (eds.) Machine Learning Proceedings 1995, pp. 331–339. Morgan Kaufmann, San Francisco (CA) (1995)

29. Lazaridou, A., Kuncoro, A., Gribovskaya, E., Agrawal, D., Liska, A., Terzi, T., Gimenez, M., de Masson d'Autume, C., Kociský, T., Ruder, S., Yogatama, D., Cao, K., Young, S., Blunsom, P.: Mind the gap: Assessing temporal generalization in neural language models. In: Neural Information Processing Systems (2021)

30. Lewis, D.D., Yang, Y., Rose, T.G., Li, F.: Rcv1: A new benchmark collection for text categorization research. J. Mach. Learn. Res. **5**, 361–397 (dec 2004)

31. Li, X., Li, Z., Sheng, J., Slamu, W.: Low-resource text classification via cross-lingual language model fine-tuning. In: Proceedings of the 19th Chinese National Conference on Computational Linguistics. pp. 994–1005. Chinese Information Processing Society of China, Haikou, China (Oct 2020)

32. Liu, J., Chang, W.C., Wu, Y., Yang, Y.: Deep learning for extreme multi-label text classification. In: Proc. 40th International ACM SIGIR Conference on Research and Development in Information Retrieval. p. 115–124. SIGIR '17, ACM, New York, NY, USA (2017)

33. Madjarov, G., Kocev, D., Gjorgjevikj, D., Džeroski, S.: An extensive experimental comparison of methods for multi-label learning. Pattern Recognition **45**(9), 3084–3104 (2012), best Papers of Iberian Conference on Pattern Recognition and Image Analysis (IbPRIA'2011)

34. Manning, C.D., Raghavan, P., Schütze, H.: Introduction to Information Retrieval. Cambridge University Press (2008)

35. Mirzakhalov, J.: Turkic Interlingua: A Case Study of Machine Translation in Low-resource Languages. Ph.D. thesis, University of South Florida (2021)

36. Mirzakhalov, J., Babu, A., Ataman, D., Kariev, S., Tyers, F., Abduraufov, O., Hajili, M., Ivanova, S., Khaytbaev, A., Laverghetta Jr, A., et al.: A large-scale study

of machine translation in turkic languages. In: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. pp. 5876–5890 (2021)

37. Mirzakhalov, J., Babu, A., Kunafin, A., Wahab, A., Moydinboyev, B., Ivanova, S., Uzokova, M., Pulatova, S., Ataman, D., Kreutzer, J., et al.: Evaluating multiway multilingual nmt in the turkic languages. In: Proceedings of the Sixth Conference on Machine Translation. pp. 518–530 (2021)

38. Moskvichev, A., Dubova, M., Menshov, S., Filchenkov, A.: Using linguistic activity in social networks to predict and interpret dark psychological traits. In: Artificial Intelligence and Natural Language: 6th Conference, AINL 2017, St. Petersburg, Russia, September 20–23, 2017, Revised Selected Papers 6. pp. 16–26. Springer (2018)

39. Nikolenko, S.I.: Topic quality metrics based on distributed word representations. In: Proc. 39th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 1029–1032 (2016)

40. Nikolenko, S.I., Koltsova, O., Koltsov, S.: Topic modelling for qualitative studies. Journal of Information Science **43**(1), 88–102 (2017)

41. Nivre, J., Zeman, D., Ginter, F., Tyers, F.: Universal Dependencies. In: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Tutorial Abstracts. Association for Computational Linguistics, Valencia, Spain (Apr 2017)

42. Oleynik, M., Kugic, A., Kasáč, Z., Kreuzthaler, M.: Evaluating shallow and deep learning strategies for the 2018 n2c2 shared task on clinical text classification. Journal of the American Medical Informatics Association **26**(11), 1247–1254 (09 2019)

43. OpenAI: Gpt-4 technical report (2023)

44. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. Journal of Machine Learning Research **12**, 2825–2830 (2011)

45. Polat, Y., Zakirov, A., Bajak, S., Mamatzhanova, Z., Bishkek, K.: Machine translation for kyrgyz proverbs—google translate vs. yandex translate-from kyrgyz into english and turkish. In: Материалы Шестой Международной конференции по компьютерной обработке тюркских языков «TurkLang-2018»(Ташкент, Узбекистан, 18–20 октября 2018 г.) (2018)

46. Read, J., Pfahringer, B., Holmes, G., Frank, E.: Classifier chains for multi-label classification. Machine Learning **85**, 333–359 (2011)

47. Reddi, S.J., Kale, S., Kumar, S.: On the convergence of adam and beyond. In: International Conference on Learning Representations (2018)

48. Reimers, N., Gurevych, I.: Sentence-bert: Sentence embeddings using siamese bert-networks. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics (11 2019)

49. Rusnachenko, N., Loukachevitch, N., Tutubalina, E.: Distant supervision for sentiment attitude extraction. In: Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019). pp. 1022–1030. INCOMA Ltd., Varna, Bulgaria (Sep 2019). https://doi.org/10.26615/978-954-452-056-4_118, https://aclanthology.org/R19-1118

50. Sadykov, T., Kochkonbayeva, B.: Model of morphological analysis of the kyrgyz language. In: Proceedings of the V International Conference on Computer Processing of Turkic Languages Turklang. vol. 2, pp. 135–154 (2017)

51. Savchenko, A., Alekseev, A., Kwon, S., Tutubalina, E., Myasnikov, E., Nikolenko, S.: Ad lingua: Text classification improves symbolism prediction in image advertisements. In: Proceedings of the 28th International Conference on Computational Linguistics. pp. 1886–1892 (2020)
52. Sennrich, R., Haddow, B., Birch, A.: Neural machine translation of rare words with subword units. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 1715–1725 (2016)
53. Shen, Z., Zhang, S.: A novel deep-learning-based model for medical text classification. In: Proceedings of the 2020 9th International Conference on Computing and Pattern Recognition. p. 267–273. ICCPR '20, Association for Computing Machinery, New York, NY, USA (2021)
54. Song, K., Tan, X., Qin, T., Lu, J., Liu, T.Y.: Mpnet: Masked and permuted pretraining for language understanding. arXiv preprint arXiv:2004.09297 (2020)
55. Szymański, P., Kajdanowicz, T.: A scikit-based Python environment for performing multi-label classification. ArXiv e-prints (Feb 2017)
56. Tang, P., Jiang, M., Xia, B.N., Pitera, J.W., Welser, J., Chawla, N.V.: Multi-label patent categorization with non-local attention-based graph convolutional network. Proc. AAAI Conference on Artificial Intelligence **34**(05), 9024–9031 (Apr 2020)
57. Toleush, A., Israilova, N., Tukeyev, U.: Development of morphological segmentation for the kyrgyz language on complete set of endings. In: Intelligent Information and Database Systems: 13th Asian Conference, ACIIDS 2021, Phuket, Thailand, April 7–10, 2021, Proceedings 13. pp. 327–339. Springer (2021)
58. Tukeyev, U., Karibayeva, A., Zhumanov, Z.h.: Morphological segmentation method for turkic language neural machine translation. Cogent Engineering **7**(1), 1856500 (2020)
59. Tutubalina, E., Nikolenko, S.: Inferring sentiment-based priors in topic models. In: Pichardo Lagunas, O., Herrera Alcántara, O., Arroyo Figueroa, G. (eds.) Advances in Artificial Intelligence and Its Applications. pp. 92–104. Springer International Publishing, Cham (2015)
60. Tutubalina, E., Nikolenko, S.I.: Constructing aspect-based sentiment lexicons with topic modeling. In: Proc. 5th International Conference on Analysis of Images, Social Networks, and Texts. pp. 208–220 (2016)
61. Tutubalina, E., Nikolenko, S.I.: Exploring convolutional neural networks and topic models for user profiling from drug reviews. Multimedia Tools and Applications **77**(4), 4791–4809 (2018)
62. Vu, H.T., Nguyen, M.T., Nguyen, V.C., Pham, M.H., Nguyen, V.Q., Nguyen, V.H.: Label-representative graph convolutional network for multi-label text classification. Applied Intelligence **53**(12), 14759–14774 (nov 2022)
63. Ward Jr, J.H.: Hierarchical grouping to optimize an objective function. Journal of the American statistical association **58**(301), 236–244 (1963)
64. Washington, J.N., Salimzianov, I., Tyers, F.M., Gökırmak, M., Ivanova, S., Kuyrukçu, O.: Free/open-source technologies for turkic languages developed in the apertium project. In: Proceedings of the International Conference on Turkic Language Processing (TURKLANG 2019) (2019)
65. Washington, J.N., Ipasov, M., Tyers, F.M.: A finite-state morphological transducer for kyrgyz. In: LREC. pp. 934–940 (2012)
66. Yao, L., Mao, C., Luo, Y.: Clinical text classification with rule-based features and knowledge-guided convolutional neural networks. CoRR **abs/1807.07425** (2018), http://arxiv.org/abs/1807.07425

67. Yin, W., Hay, J., Roth, D.: Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 3914–3923 (2019)

68. Yiner, Z., Kurt, A., Kulamshaev, K., Zafer, H.R.: Kyrgyz orthography and morphotactics with implementation in nuve. In: Proceedings of International Conference on Engineering and Natural Sciences. pp. 1650–1658 (2016)

69. Zhang, M.L., Zhou, Z.H.: Multilabel neural networks with applications to functional genomics and text categorization. IEEE Transactions on Knowledge and Data Engineering **18**(10), 1338–1351 (2006). https://doi.org/10.1109/TKDE.2006.162

70. Zhang, M.L., Zhou, Z.H.: Ml-knn: A lazy learning approach to multi-label learning. Pattern recognition **40**(7), 2038–2048 (2007)

71. Zhang, Y., Surendran, A.C., Platt, J.C., Narasimhan, M.: Learning from multi-topic web documents for contextual advertisement. In: Proceedings of the 14th ACM SIGKDD international conference on knowledge discovery and data mining. pp. 1051–1059 (2008)

72. Zhu, H., Lei, L.: The research trends of text classification studies (2000–2020): A bibliometric analysis. SAGE Open **12**(2), 21582440221089963 (2022)

73. Арыкоглу, Экрем: Словарный проект современного кыргызского языка. In: Общество, язык и культура XXI века. pp. 85–91 (2021)

74. Бакасова, П. С., Исраилова, Н. А.: Алгоритм образования словоформ для автоматизации процедуры пополнения базы данных словаря. Известия Кыргызского государственного технического университета им. И. Раззакова (2), 23–27 (2016)

75. Исраилова, Н. А.: Алгоритм морфологического анализа и синтеза в трансляторе. Современные проблемы механики (28), 11–19 (2017)

76. Исраилова, Н. А., Бакасова, П. С.: Онтологические модели морфологических правил кыргызского языка. In: Седьмая Международная конференция по компьютерной обработке тюркских языков «TurkLang 2019»: Труды конференции (Симферополь, Крым, Россия, 3–5 октября 2019 г.) (2019)

77. Карабаева, С. Ж. : Использование грамматических правил в Прологе. Вестник Бишкекского гуманитарного университета (2), 231–233 (2011)

78. Кочконбаева, Б. О.: Табигый тилдеги текстттерди орус тилинен кыргыз тилине машиналык которууда сездерду анализдве^ н алгоритмин тузуу. Известия Кыргызского государственного технического университета им. И. Раззакова (2), 52–54 (2016)

79. Кочконбаева, Б. О., Эгембердиева, Ж.С.: Modeling of morphological analysis and synthesis of word forms of the natural language. Бюллетень науки и практики **6**(9), 435–439 (2020)

80. Момуналиев, К. З.: Парсирование и аннотирование турецко-кыргызского словаря. Известия Кыргызского государственного технического университета им. И. Раззакова (2), 68–81 (2016)

81. Мусаев, С. Ж., Карабаева, С. Ж., Иманалиева, А.И.: Проблемы и перспективы развития компьютерной лингвистики в Кыргызстане. In: PROCEEDINGS Of the I International Conference on Computer processing of Turkic Languages (TurkLang-2013). pp. 34–37 (2013)

82. Садыков, Т, Кочконбаева, Б: Об оптимизации алгоритма морфологического анализа. In: Сборник материалов Шестой Международной конференции по

компьютерной обработке тюркских языков «TurkLang-2018»(Ташкент, Узбекистан, 18–20 октября 2018 г.) (2018)

83. Садыков, Т., Шаршембаев, Б.: «Манас» эпосунун улттук корпусун түзүү жөнүндө. In: Компьютерная обработка тюркских языков. Первая меж дународная конференция: Труды. pp. 148–154. No. 6, Астана: ЕНУ им. ЛН Гумилева. (2013)

84. Шарипбай, А. А., Ергеш, Б. Ж., Елибаева, Г. К., Жеткенбай, Л., Исраилова, Н., Бакасова, П.: Сравнение онтологических моделей существительных казахского и кыргызского языков. In: Материалы Шестой Международной конференции по компьютерной обработке тюркских языков «TurkLang-2018»(Ташкент, Узбекистан, 18–20 октября 2018 г.) (2018)