**UBMK'24**

**Bildiriler Kitabı**
**Proceedings**

Editor Eşref ADALI

# 9. Uluslararası Bilgisayar Bilimleri ve Mühendisliği Konferansı

# 9th International Conference on Computer Science and Engineering

26-27-28 Ekim (October) 2024 Antalya - Türkiye

# 9. Uluslararası Bilgisayar Bilimleri ve Mühendisliği Konferansı (UBMK'2024)

# 9th International Conference on Computer Science and Engineering

26-28 Ekim 2024 Akdeniz Üniversitesi Antalya Türkiye
26-28 October 2024 Akdeniz University Antalya Türkiye

| | | | | |
|---|---|---|---|---|
| | 2150 | Business Process Management Anomaly Detection through Semantic Embedding-Integrated Graph Neural Networks | Teoman Berkay Ayaz | 568 - 573 |
| | | | Ege Gülce<br>Stanley Hsu<br>Alper Özcan<br>Akhan Akbulut | |
| | 2155 | Multi-Aspect Anomaly Detection with Graph Neural Networks and Kolmogorov-Arnold Networks in Business Process Management | Teoman Berkay Ayaz | 574 - 579 |
| | | | Ege Gülce<br>Stanley Hsu<br>Alper Özcan<br>Akhan Akbulut | |
| | 2156 | Comparison of a Deep Learning and a Hybrid Model for Classification of an Unbalanced Urgent Cases Dataset for Human Faces | Faruk Özgür | 580 - 585 |
| | | | Neslihan Arıkan<br>Özge Öztimur Karadağ | |
| | 2171 | Estimating the Manufacturing Cost of a Metal Part from Textual and Geometric Features | Talha Rehman Abid | 586 - 591 |
| | | | Mert Daloğlu<br>Cem Yıldız<br>Ali Erman Erten<br>Kamer Kaya | |
| | 2176 | Object Detection in Hyperspectral Images with Unsupervised Domain Adaptation | Sinem Aybüke Şakacı<br>Alp Ertürk<br>Erchan Aptoula | 592 - 596 |
| | 2182 | Machine Learning Approaches to Predict Thyroid Cancer Recurrence: A Comparative Study | Candide Ozturk | 597 - 602 |
| | | | Ozgur Sagir<br>Ulas Vural | |
| | 2197 | Yangın Söndürme Süreçlerinde Su Tüketiminin Makine Öğrenmesi ile Tahmini Prediction of Water Consumption in Fire Extinguishing Processes Using Machine Learning Approaches | Emin Ölmez | 603 - 606 |
| | | | Ahsen Usta | |
| | | | Orhan Akbulut | |
| | 2202 | Synthetic Vibration Data Generation and Fault Classification in CNC Machines Using Transformer GANs and ConvLSTM Networks | Özlem Erbay | 607 - 612 |
| | | | Batıray Erbay | |
| IoT | 1862 | Design and Implementation of a Management System for Wireless Electronic Combination Locks | Batuhan Kol | 613 - 618 |
| | | | Metin Bilgin | |
| | 1887 | Solar IoT: Monitoring the Orientation and Electrical Parameters of the Solar Panel | Hakan Dalkılıç<br>Oğuz Gora | 619 - 624 |
| | 1973 | Multifunctional Smart Spoon for Parkinson's: Stability Enhancement and Diagnostic Tools | Divyansh Singhal | 625 - 629 |
| | | | Sasi Snigdha Yadavalli<br>Nupur Patil<br>Siddharth Chauhan<br>Madhav Rao | |
| | 1988 | An AI-Assisted Autonomous IoRT Agent for Smart Spaces | Yakup Kayataş<br>Sanem Kabadayı | 630 - 635 |
| | 1991 | Advanced Detection and Prevention of Sinkhole Attacks in 6TiSCH Networks | Burak Aydın<br>Hakan Aydın<br>Sedat Görmüş | 636 - 641 |
| | 2000 | Comparison of Different Weather Data Acquisition Methods | Emre Evcin<br>Yusuf Murat Erten | 642 - 647 |
| | 2057 | Employing Digital Twin to Forest Fire Management Systems | Bugra Aydin<br>Sema Fatma Oktug | 648 - 653 |
| | 2085 | Distributed Key Value Store for IoT Edge Devices | Burak Aslantaş<br>Elif Nurdan Pektaş<br>Şebnem Baydere | 654 - 659 |
| EMB | 2096 | Soft Error Reliability Assessment of TinyML Algorithms on STM32 Microcontroller | Ahmet Selim Karakuş<br>Osman Buğra Göktaş<br>Sadık Akgedik<br>Sanem Arslan | 660 - 664 |
| HUM | 1999 | Development of An Algorithm for Converting Json Formats to Xml by Forming its File Data Structure | Aigul Mukhitova | 665 - 670 |
| | | | Aigerim Yerimbetova<br>Vladimir Barakhnin | |

| | | | | |
|---|---|---|---|---|
| | | ChatGPT Supported Interpretation on Anomaly Detection in Retail Data and Exchange Rate Relationship | Hatice Nizam-Özoğur | |
| | 2033 | Text to SQL Transformation Using LLM: a Comparative Research of T5, Seq2Seq, and SQLNet Models | Zhazira Shaikhiyeva | 967 - 972 |
| | | | Madina Mansurova Gulshat Amirkhanova | |
| | 2076 | Sağlık Sigortası Sahiplerinin Davranışsal Analizi ve Kümelenmesi Clustering and Behavioral Analysis of Health Insurance Owners | Omer Sezer Koyuncu Seçil Arslan | 973 - 978 |
| | 2087 | On symbolic Prediction of Time Series for Predictive Maintenance Based on SAX-LSTM | Aykut Güler | 979 - 983 |
| | | | Tuğçe Ballı E. Fatih Yetkin | |
| | 2135 | Profiling Driver Behaviors Using AI-Based Methods and Deep Learning Techniques for Improving Road Safety: A Comparative Study of Algorithms | Volkan Oban | 984 - 989 |
| | | | Mustafa Kaya Güzide Safi İrem Nur Çimen Tubanur Çatak Bulut Karadağ Gökhan Gümüş Aslıhan Çandır Fatih Alagöz | |
| IR | 1896 | ReRag: A New Architecture for Reducing the Hallucination by Retrieval-Augmented Generation | Robin Koç | 990 -994 |
| | | | Mustafa Kağan Gürkan Fatoş T. Yarman Vural | |
| | 1941 | Enhancing Object Detection in Aerial Images Using Transformer-Based Super-Resolution | Aslan Ahmet Haykır | 995 - 1000 |
| | | | İlkay Öksüz | |
| NET | 1985 | Proof of Concept Implementation for RSVP TSN Control Plane | Necip Gozuacik | 1001 - 1004 |
| | 2100 | Integrating Blockchain and SDN for Centrality-Aware Virtual Multicast Tree Embedding | Furkan Ayaz | |
| | | | Evrim Guler Murat Karakus Davut Hanbay | 1005 - 1010 |
| | 1969 | QoS Aware Routing Approaches in Software Defined Smart Grids | Sedef Demirci | 1011 - 1016 |
| | 2008 | Deep Reinforcement Learning Routing in Mobile Networks | Arif Burak Dikmen Hasari Çelebi | 1017 - 1022 |
| RBOT | 1942 | Endüstriyel Robotik Sistemlerin Güvenlik Doğrulaması Safety Verification of Industrial Robotic Systems | Fatih Furkan Arslan Metin Özkan | 1023 - 1028 |
| | 2077 | EKF Based Localization: Integrating IMU and LiDAR Data in the Hilti SLAM Challenge | Behice Bakır Havvanur Bozömeroğlu Ebu Yusuf Güven | 1029 - 1034 |
| SING | 1965 | Communication (Educational) Kit (HaKi) | Murat Sever Utku Bilgin | 1035 - 1038 |
| | 2089 | Manyetik Parçacık Görüntülemede Sistem Matrisi için Farklı Dalgacık Dönüşümlerinin Seyreklik Seviyesi Karşılaştırması Sparsity Level Comparison of Different Wavelet Transforms for the System Matrix in Magnetic Particle Imaging | Vildan Atalay Aydın | 1039 - 1043 |
| | 2097 | Sparse Channel Estimation For M-QAM-Based Underwater Acoustic Communication Systems | Mhd Tahssin Altabbaa | 1044 - 1048 |
| | | | Berkay Tekat Emin Tarik Iseri | |
| OTH | 1858 | The 80/20 Principle in Morphemics-Morphology in the Educational Corpus of the Uzbek Language | Shahlo Khamroeva | 1049 - 1052 |
| | | | Bakhtiyor Mengliyev Muyassar Kholova | |
| | 1904 | Gamification as a Tool for Personalized Learning in Inclusive Education | Dilaram Baumuratova Tamara Zhukabayeva Mira Rakhimzhanova | 1053 - 1058 |
| | 1918 | A Metaheuristic Algorithm for the Fixed Charge Transportation Problem | Nermin Kartli | 1059 - 1062 |
| | 2027 | Eğitimde Sürükleyici Teknolojilerin Kullanılması Fırsatlar ve Beklentiler | Atamuratov Rasuljon Kadırjanovich Majıdova Gulhayo Abdırazzoq qızı Bayjonov Furqat Baxramovıch Ongarov Mansurbek Bayrambekovıch | 1063 - 1068 |

| | | | | |
|---|---|---|---|---|
| | | | Saydullayev Zafar Erkınovıch | |
| | 2103 | Bilgisayar Mühendisliği Öğrencilerinin Perspektifinden Bilişim Hukukunun Güncel Sorunları ve Çözüm Önerileri<br>Current Challenges and Solution Proposals in IT Law from the Perspective of Computer Engineering Students | Sevda Bora Çınar | 1069 - 1075 |
| | 2200 | A Comparison of shcU-Net Based GAN and U-net Based GAN in Adult Dental Segmentation | Gürdal Altundağ<br><br>Hakan Öcal | 1075 - 1080 |
| | 1932 | Leveraging Quantum Computing and Optimization to Estimate Financial Crashes in Small and Medium-Sized Enterprises | Ege Dincer<br><br>Berkay Coskuner<br>Ege Bilaloglu<br>Bilge Koroglu | 1081 - 1086 |
| SW | 1859 | Investigating The Adoption of International Software Quality Standards in Turkey: A Comprehensive Analysis | Sevgi Koyuncu Tunç | 1087  - 1093 |
| | 1886 | Development of the Functional Structure of the Science and Education Information System | Dauletov Adilbek Yusupbayevich<br><br>Matyakubova Noila Shakirjanovna | 1094 - 1098 |
| | 1892 | React ve Preact Javascript Çerçevelerinde Karşılaştırmalı Analiz<br>Comparative Analysis on React and Preact Javascript Frameworks | Muhammed Furkan Uygur<br>Nesibe Yalçın | 1099 - 1104 |
| | 1917 | CAGE: A Tool for Code Assessment and Grading | Ümit Kanoğlu<br>Oğuz Kerem Yıldız<br>Hasan Sözer<br>Olcay Taner Yıldız | 1115 - 1110 |
| | 1957 | Extracting Driving Styles from Automotive Sensor Data to Develop Personas | M. Cagri Kaya<br>Tayssir Bouraffa<br>Krzysztof Wnuk | 1111 - 1114 |
| | 1962 | Lojistik Sipariş Dağıtım Entegrasyonu Sürecinde Sipariş Geri Çağırma Süreci Tasarımı ve Yazılım Geliştirmesi<br>Design and Software Development of The Order Recall Progress in The Logistics Order Distribution Integration Process | İklim Barman<br><br>Ersin Şengül | 1115 - 1120 |
| | 2009 | The Dimension of Green Coding in Software Quality Control Processes | Volkan Abur | 1121 - 1126 |
| | 2055 | Are We Asking the Right Questions to ChatGPT for Learning Software Design Patterns? | Çağdaş Evren Gerede | 1127 - 1132 |
| | 2060 | Optimizing LLVM IR: Transforming Multiplication to Addition for Enhanced Execution Efficiency | Huseyin Karacalı<br><br>Efecan Cebel<br>Nevzat Donum | 1133 - 1138 |
| | 2080 | Estimation of Software Integration Test Duration via UML Statecharts | Fehim Göler<br>Tolga Ovatman | 1139 - 1144 |
| | 2093 | DIA4M: A Tool to Streamline DevOps Processes of Distributed Cloud-Native Systems | Eren Tarak<br>H. Hakan Kilinc | 1145 - 1150 |
| | 2111 | Software Industry Perception of Academic Collaboration | Deniz Akdur | 1151 - 1156 |
| | 2139 | Görüntü İşlemeyle Doğrulamalı Robotik Test  Otomasyon Kullanımı: POS Cihazları Üzerine Uygulama | Miraç Emektar<br><br>Harun Kadıoğlu<br>Ahmet Efendioğlu<br>Fatih Mehmet Harmancı | 1157 - 1161 |
| | 2141 | VoIP Sistemlerinde Zihin Haritası Tabanlı Test  Stratejiler : SIP Pbx Ürünü Üzerine Bir İnceleme<br><br>Mind Map-Based Testing Strategies  in VoIP Systems: A Case Study on SIP  Pbx Products | Miraç Emektar<br><br>Furkan Günaydın<br><br>Fatih Mehmet Harmancı | 1162 - 1167 |
| | 2173 | A Robust Microservices Framework for Indoor Tracking System Development | Gafur Hayytbayev<br>Kerem Küçük<br>Mahmut Çavur | 1168 - 1172 |
| DM | 1927 | Unsupervised Pattern Extraction of Time Series Data for Energy Disaggregation | Şirin Azazi Deveci<br>Melih Günay | 1173 - 1178 |
| | 1944 | Topic Modeling Enhanced Tripartite Graph for Recommendation using Metapaths | Yaren Yılmaz<br>Irem İşlek<br>Şule Gündüz Öğüdücü | 1179 - 1184 |
| | 1948 | Community Detection on Software Library Dependency Graphs using Graph Neural Networks | Şevket Umut Çakır | 1185 - 1190 |

| | | | | |
|---|---|---|---|---|
| | | | Mehmet Ali Osman Atik<br>Ümit Deniz Uluşar | |
| | 2190 | Enhancing Mesh and Point Cloud Similarity Detection through Geometric Features and ICP | Talha Rehman Abid | 1191 - 1196 |
| | | | Mehtap Öklü<br>Cem Yıldız<br>Ali Erman Erten<br>Kamer Kaya | |
| | 2214 | Comparative Analysis and Practical Implementation of Machine Learning Algorithms for Phishing Website Detection | Samad Najjar-Ghabel | 1197 - 1202 |
| | | | Shamim Yousefi<br>Payam Habibi | |
| | 2215 | A Technical Analysis and Practical Implementation of Machine Learning Algorithms for Predicting Survival in Breast Cancer Patients | Shamim Yousefi | 1203 - 1208 |
| | | | Samad Najjar-Ghabel<br>Hamidreza Shafaei | |
| BIG | 1881 | Comparison Between Time Series and Relational Databases | Alpar Türkoğlu<br>Onurcan Ersen<br>Ibrahim Onuralp Yiğit<br>Dincer Unal<br>Hatice Golcuk | 1209 - 1212 |
| | 1930 | A Performance Evaluation Study on a Data Analytics Platform for Emergency Calls | Engin Yakar<br>H. Hakan Kilinc | 1213 - 1218 |
| | 2079 | Adaptive Composite Market Volatility Index (CMVI) for Enhanced Stock Price Forecasting | Rabia Çevik | 1219 - 1223 |
| | | | Uğur Barış Özyürek<br>Ali Kanal<br>Vael Kokach<br>Büşra Kocaçınar<br>Oznur Şengel<br>Fatma Patlar Akbulut | |
| | 2142 | Hybrid Deep Learning Framework for Stock Price Prediction Incorporating Technical and Macroeconomic Indicators | Ali Can Turan | 1224 - 1228 |
| | | | Vael Kokach<br>Büşra Kocaçınar<br>Oznur Şengel<br>Fatma Patlar Akbulut | |
| | 2125 | Emotion-Aware Multimodal Biometric Identification by using Biosignals | Yekta Said Can<br>Beyzanur Bektan<br>Fatih Alagöz | 1229 - 1235 |
| | 1854 | Özbekçe-Türkçe Otomatik Çeviri Yazılımı için Deyimlerin Veritabanını Teşkil Etmede Karşılaşılan Güçlükler Automatic Translation Software Difficulties in Organizing the Database of Idioms for Uzbek and Turkish | Manzura Abjalova | 1236 - 1240 |
| | | | Umida Raşidova<br>Eşref Adalı | |
| | 2028 | Reversible Steganographic System for the Transmission of Personal Medical Data | Elmira Daiyrbayeva<br>Ekaterina  Merzlyakova<br>Aigerim Yerimbetova<br>Aigul Mukhitova | 1241 - 1246 |

# Araneum Uzbecicum: A Gigaword Web-Crawled Uzbek Corpus

Vladimír Benko
*L'. Štúr Institute of Linguistics*
*Slovak Academy of Sciences*
Batislava, Slovkia
vladimir.benko@juls.savba.sk

Radovan Garabík
*L'. Štúr Institute of Linguistics  Slovak Academy*
*of Sciences* Batislava, Slovkia
radovan.garabik@kassiopeia.juls.savba.sk

Shahlo Khamroyeva
*Tashkent State University of Uzbek*
*language and literature*
Tashkent, Uzbekistan
shaxlo.xamrayeva@navoiy-uni.uz

*Abstract*— **We want to introduce a Project of creation of a web-crawled Uzbek corpus. The main steps related to data capture, pre-processing, tokenization, deduplication, PoS tagging and lemmatization are discussed, and some query examples are shown. This article analyzes the problems in the process of creating and lemmatization the corpus of Araneum Uzbecicum. We also show word embeddings created from the Uzbek data. The Word embeddings function serves as a useful tool for researching the semantic relations in the Uzbek language, word context and valence. Also, the paradigm created by the Word embeddings function helps to perfect dictionaries.**

*Keywords—Uzbek language, web-crawled corpus, PoS tagging and lemmatization, word embeddings*

## I. INTRODUCTION

After Tatar and Kazakh, Uzbek is the third Turkic language in our corpus. Spoken by approx. 33 million of L1 or L2 speakers in Uzbekistan and several of its neighboring countries, and also being "sufficiently represented" in the Internet, the language is a suitable target for creation of a web-crawled corpus.

In our paper, we would like to introduce a Project aimed at creation of a Gigaword Uzbek corpus that would contain texts in both scripts currently used in Uzbekistan, i.e., modified Latin and modified Cyrillic, respectively. The corpus data is PoS tagged and lemmatized by *Apertium*, processed by the *NoSketch Engine* corpus manager and made available at our web-corpora portal[1].

## II. RELATED WORK

In recent years, building web corpora has become a well-established part of linguistics. The Web, teeming as it is with language data, of all manner of varieties and languages, in vast quantities and freely available, is a fabulous linguists' playground. This special issue of Computational Linguistics explores ways in which this dream is being explored.

The use of web corpora in linguistic research has become widespread, so the problem of using web corpora in linguistic research has been studied in depth.

The problems of creating a web corpus using the BootCaT tool for low-resource languages were studied and clear results were obtained.

Several corpora of Turkish languages were developed on the Sketch Engine platform.

"uzb_community_2017" is an Uzbek community corpus based on material from 2017. It contains 663,119 sentences and 9,256,001 tokens [2].

The initial efforts to develop the Uzbek language corpus commenced in Uzbekistan in 2018. While there have been some theoretical articles on the construction of the Uzbek language corpus, the practical developments on the corpus have become more prevalent by 2018. A bibliometric analysis based on Scopus was conducted on corpus linguistics for the period spanning 2017 to 2021 [41]. There is some work about creating a morphological and syntactic tagged corpus for the Uzbek language. Several studies have been conducted on the processing of corpus texts in the Uzbek language [44], [45].

Additionally, several other corpora have been created for the Uzbek language. In particular, the "Educational corpus of the Uzbek language"[3] within the project of the Tashkent State University of Uzbek Language and Literature, the "Uzbek language corpus"[4]; based on the project of the National University of Uzbekistan, the "National corpus of the Uzbek language". "Alisher Navoi authorship corpus"[5] and the "Diachronic corpus of Uzbek-English newspapers"[6] were developed by independent researchers, are available for free use. Software for the national corpus of the Uzbek language was developed.

Furthermore, it is worth noting the contributions of researchers who have undertaken studies on specific types of corpora. These include articles on the development of the Uzbek author's corpus, the Uzbek-English parallel corpus, Uzbek-Russian parallel corpus and the diachronic corpus of the Uzbek language. There is research in the field of author lexicography and the relationship between author corpora (Anonymous). An algorithm for creating a parallel corpus of the Uzbek and Russian languages was. The theoretical foundations of the Uzbek-English parallel corpus have been developed [42], [43]. A scientific study of the open Uzbek speech corpus and preliminary experiments on speech recognition were conducted.

## III. DATA FOR THE NEW CORPUS

During the last decade, the methodology and tools for compilation web-crawled corpora has been effectively standardized. In our Project, we were using a collection of tools usually referred to as "Brno Pipeline" available under a

---

[1] http://aranea.juls.savba.sk/guest/index.html
[2] https://www.sketchengine.eu/uzwac-uzbek-corpus/
[3] https://uzschoolcorpara.uz/

[4] https://uzbekcorpus.uz/
[5] https://alishernavoiykorpusi.uz/
[6] https://mediatextcorpus.com/

FLOSS license[7] . Its main component is SpiderLing[8], a web crawler optimized for downloading textual data. In its latest Version 2.2, it represents a mature, effective and stable tool that is able to get (within our technical infrastructure) as much as two Gigaword of plain text within 24 hours of crawling. The tool also includes a character encoding identification utility; a language detection module that is using samples of texts for the languages to be recognized provided by the user to build the respective language models using character trigram frequencies; and a boilerplate identification and removal tool .. All the processes work smoothly on the fly, outputting a reasonably clean plain text with only light XML markup encoding some metadata at the document and paragraph level.

To initialize the process, SpiderLing needs a list of seed URLs that are used in the first round of crawling. In our case, the seed list of approx. 2,000 URLs has been obtained by means of the WebBootCat functionality accessible at the Sketch Engine Portal.

The actual crawling was performed during two separate sessions, in December 2020, and March 2024, respectively. The seed URLs for the second one was obtained from the data already downloaded during the first session – the last 100,000 URLs were used, expecting that some of the addresses might meanwhile have changed the contents of the respective web pages, or could be removed during deduplication.

## IV. PRE-TOKENIZATION PROCESSING

Though the language identification procedure of SpiderLing works reliably for most languages, it is often not able to discriminate between very close languages using the same script. We therefore perform a secondary language filtering based on character frequencies, making use of the fact that some characters are unique for a particular language. This was also the case of Uzbek that is using the x, к, ў, and ғ in the Cyrillic script, and o', g' and ' characters or character combinations in the Latin script. The respective thresholds for filtration were established experimentally.

Other filters detect potential mojibake encoding problems appearing as "artifacts" such as *subâ€TMektlarini* or *oÃf¢€TMtgan* in the texts – documents with too many suspicious character sequences are identified by means of a series of regular expressions and removed.

## V. TOKENIZATION

For various reasons, we perform tokenization as a separate processing step, even in situations when the tagger would prefer to do it as a part of the tagging process. Here again, one of the "Brno" tools is used – Unitok ., the universal tokenizer is a Python program using a custom-modified parameter file containing regular expressions for language-specific period-final abbreviations, ordinal numbers, etc. The initial list of Uzbek abbreviations (both in Latin and Cyrillic) presently contains just approx. two dozen items – we hope, however, to be able to expand it after analyzing the processed corpus.

## VI. DATA NORMALIZATION

Now it is necessary to mention one of the peculiarities of the Uzbek language – the simultaneous use of two different scripts, with the speakers being proficient in both of them. This situation is similar to that of the Serbian language, but unlike in Serbia (and Bosnia and Herzegovina), where the digraphia is a result of centuries long development and any of the scripts are widely accepted in almost all situations, official and unofficial, , the de facto digraphia in Uzbekistan is a result of partially unsuccessful transition from Cyrillic to Latin script in early 1990s. According to our observation, the amount of texts in Internet in either of the scripts is about equal, with many websites, including those belonging to government organizations, private companies, or even the Uzbek Wikipedia, provide two versions of their web pages.

We have decided, on one hand, to retain the information on the respective script in our corpus, as we consider it to be an interesting sociolinguistic parameter, and on the other hand to provide transliterated version of all texts, so that users could easily query the corpus. Moreover, this is also necessary as the PoS tagger used expects its input in Latin script.

The other issue is the presence of the "modified letter turned comma" (U+02BB) and "modified letter apostrophe" (U+02BC) characters in the Latin script. Since these characters are absent from most keyboard layouts, they are often replaced by similar ones, such as "left single quotation mark" (U+0218), "apostrophe" (U+0027), or "right single quotation mark" (U+0219), or oven other similarly looking characters, such as an acute or grave accent. As the PoS tagger expects the text to contain canonical forms of these characters only, all these irregularities have to be normalized before further processing.

## VII. DEDUPLICATION

Web-crawled texts are known to contain a large proportion of duplicate contents at different levels. We have adopted a two-step approach: (partial) duplicates at the document level are removed and not passed to further processing, while paragraph and sentence level deduplication is provided, if required, at the very end of the processing only.

For the document-level deduplication we use a tool called Onion working on the principle of comparing word n-grams from the document processed with those encountered in all previous documents. A document is considered (partially) duplicate if it contains certain proportion of n-grams already seen. Both values, i.e., the length of n-grams to compare and the similarity level is set by the user – we work with 5-grams and 90% threshold as discussed in [45].

## VIII. POS TAGGING AND LEMMATIZATION

Uzbek unfortunately does not belong to languages with a language model available for the "traditional" multilingual taggers, such as TreeTagger], or at least a lemmatizer like CSTlemma, though we hope the situation might improve in the future.

The baseline morpho-syntactic annotation has been performed by Apertium, that is in fact machine translation tool, yet its morphological module can be used separately for morpho-syntactic annotation. Its main deficiency, for most languages where a language model is available, is a rather small size of the respective morphological lexicon and a limited disambiguation procedure. It is also not able to guess lemmas and/or PoS tags for the out-of-dictionary (OOV) lexical items.

---

[7] https://corpus.tools/

[8] https://corpus.tools/wiki/SpiderLing

Another problem of using Apertium as a tagger is its "attitude" to data already tokenized. Being a machine translation tool, its lexicon can also contain multi-word expressions than are to be translated as whole units. If such a multi-word unit is encountered in the text, the system decides that a new token to be added, i.e., spoiling the original tokenization that is expected by the overall processing, e.g., to be able to merge the Cyrillic and Latin data. As the Apertiun morphological lexicon is available in a readable form, we have decided to delete all multi-word expressions from it during this initial phase of our Project.

The respective processing steps of tokenization, normalization, PoS tagging and lemmatization are shown in Fig. 1 to 4.

```
<s>
Янги
стратегияни
август
ойига
қадар
Тошкент
шаҳри
миқёсида
тасдиқланиши
режалаштирилмоқда
<g/>
.
</s>
```

Fig. 1. Tokenized sentence in original Cyrillic Script

```
<>
Yangi
strategiyani
avgust
```

```
oyiga
qadar
Toshkent
shahri
miqyosida
tasdiqlanishi
rejalashtirilmoqda
<>

.
<>
```

Fig. 2. Normalized (Romanized) and filtered sentence to be sent as input for *Apertium*

```
<>
^Yangi/yangi<n><nom>$
^strategiyani/strategiya<n><acc>$
^avgust/*avgust$
^oyiga/oy<n><px3sp><dat>$
^qadar/qadar<post>$
^Toshkent/Toshkent<np><top><nom>$
^shahri/shahr<n><px3sp><nom>$
^miqesida/ miqyos<n><px3sp><loc>$
^tasdiqlanishi/tasdiqlan<v><iv>
    <ger_abst><px3sp><nom>$
^rejalashtirilmoqda/rejalashtir<v><tv>
    <pass><prog_inf><p3><pl>$
<>
^./.<sent>$
<>
```

Fig. 3. Sentence tagged by *Apertium* in its internal format

TABLE I. FULLY PoS TAGGED AND LEMMATIZED SENTENCE

| Word | nword | lemma | atag | tag | ztag |
|------|-------|-------|------|-----|------|
| <s> | | | | | |
| Янги | Yangi | yangi | Nn | n:nom | 1 |
| стратегияни | strategiyani | strategiya | Nn | n:acc | 1 |
| август | avgust | avgust | Yy | * | 0 |
| ойига | oyiga | oy | Nn | n:px3sp:dat | 1 |
| қадар | qadar | qadar | Pp | post | 1 |
| Тошкент | Toshkent | Toshkent | Nn | np:top:nom | 1 |
| шаҳри | shahri | shahr | Nn | n:px3sp:nom | 1 |
| миқёсида | miqyosida | miqyos | Nn | n:px3sp:loc | 1 |
| тасдиқланиши | tasdiqlanishi | tasdiqlan | Vb | v:iv:ger_abst:px3sp:nom | 1 |
| режалаштирилмоқда | rejalashtirilmoqda | rejalashtir | Vb | v:tv:pass:prog_inf:p3:pl | 1 |
| <g/> | | | | | |
| . | . | . | Zz | sent | 1 |
| </s> | | | | | |



Fig. 4. Querying the Corpus in the "Simple Query" mode

The sentence chosen as an example happened to be in Cyrillic script (Fig. 1), it needed to be converted to Latin script with the XML tags "rasped" so that Apertium would to try to tag them as well (Fig. 2). The Apertium output (Fig. 3) has to be parsed, converted to column format with the original word forms, and XML tags restored to get the format expected by the NoSketch Engine corpus manager (Table I).

The word attribute contains the word form it the original script, nword is its normalized and/or Romanized value, lemma and tag are values assigned by the tagger. Two extra attributes result from the post-PoS processing: atag is the part of speech information expressed in the Araneum Universal Tagset (Anonymous), and ztag displays the result of the morphological lexicon lookup, with the "1" value meaning success, and "0" out of vocabulary (OOV) lexical item. All these corpus attributes are available for querying to corpus users. Moreover, as word, nword, and lemma are present in the "Simple query" CQL definition, user can query either in Latin or Cyrillic script, where queries in Cyrillic script will find Cyrillic occurrences only, while those in Latin script will find occurrences in both scripts, as shown in Fig. 4.

## IX. QUERYING THE CORPUS

A 125 Megatoken sample of the corpus has already been published at our "Sandbox" Corpora Portal providing the users with full NoSketch Engine functionality. The Fig. 5, for example, shows the most salient collocates of the word "Toshkent" within the "–4 to +5" window sorted by the logDice value.

| | Cooccurrence count | Candidate count | MI | logDice |
|---|---|---|---|---|
| shahr | 6,010 | 41,068 | 8.957 | 11.304 |
| viloyat | 5,825 | 70,148 | 8.139 | 10.801 |
| shahar | 4,974 | 69,837 | 7.918 | 10.577 |
| universitet | 2,610 | 36,756 | 7.914 | 10.183 |
| institut | 1,884 | 23,981 | 8.059 | 9.988 |
| tuman | 2,902 | 63,610 | 7.275 | 9.887 |
| hokimlik | 829 | 7,715 | 8.511 | 9.253 |
| Toshkent | 1,332 | 36,754 | 6.943 | 9.212 |
| boshqarma | 808 | 15,752 | 7.444 | 8.976 |
| IIBB | 568 | 785 | 11.263 | 8.951 |
| davlat | 2,719 | 151,271 | 5.932 | 8.887 |
| hokim | 696 | 14,049 | 7.394 | 8.808 |
| shahri | 539 | 2,987 | 9.259 | 8.794 |
| Samarqand | 649 | 13,740 | 7.325 | 8.716 |
| kengash | 685 | 24,419 | 6.574 | 8.518 |
| filial | 494 | 9,627 | 7.445 | 8.445 |
| akademiya | 461 | 10,107 | 7.275 | 8.331 |
| xalqaro | 733 | 40,631 | 5.937 | 8.277 |
| nom | 1,027 | 71,913 | 5.600 | 8.274 |
| O'zbekiston | 1,478 | 122,906 | 5.352 | 8.244 |
| kun | 2,723 | 269,653 | 5.100 | 8.185 |

Fig. 5.   Most salient collocates of "Toshkent".

## X. WORD EMBEDDINGS

Word embeddings is a method that assigns multidimensional vectors to individual words. Pioneered by the word2vec, it soon became an indispensable part of many NLP applications. Although nowadays mostly replaced by contextual embeddings, word embeddings are still used for

many purposes (e.g in lexicography), because they capture semantic and grammatical properties of individual words.

We provide vector embeddings models for most of the languages of the Aranea family of corpora, the newest addition is the set of Uzbek embeddings trained on the Araneum Uzbecicum corpus with the GenSIM framework. The models use a context window 7 words wide and the skip-gram training algorithm, with dimensionality 200.

The model is accessible through a user-friendly web interface[9] that can be used to access and query the models and demonstrate several common ways of using the models. The interface is built upon an assumption that the vectors reflect semantic properties of the words and the position (and distance) of the vectors in the embedding space corresponds to the position in an abstract "semantic space" of the language. The interface is geared towards lexicographers and linguists without experience in NLP. The models use automatic detection of bigrams, these are treated as regular words (i.e. a vector is assigned to them) and can be used in queries, with the underscore character _ (U+005F LOW LINE) joining the components of the bigram.

The basic information provided is the list of words corresponding to vectors with small angular distance from the given word. Instead of cosine similarity, we define "semantic difference" as $\sqrt{(1-\cos^2 \theta)} = \sin \theta$, being close to zero for near-synonyms and close to one for unrelated words. We found out this better aligns with the expectation of our users. An example of the list of near synonyms is at Fig. X2. Then there is a visualization of the embedding vectors, we are using ISOMAP dimensionality reduction to display the results as a 2-dimensional, 3-dimensional or (for those users able to mentally imagine four-dimensional Euclidean space) 4-dimensional map, presented as a perspective projection into 2-dimensional picture, as well as two different word cloud representations. There is also a link to a raw dimensionality reduced data in Gnuplot forma, giving the users ability to pan, zoom and rotate the graphs, or use the data in other statistical software.
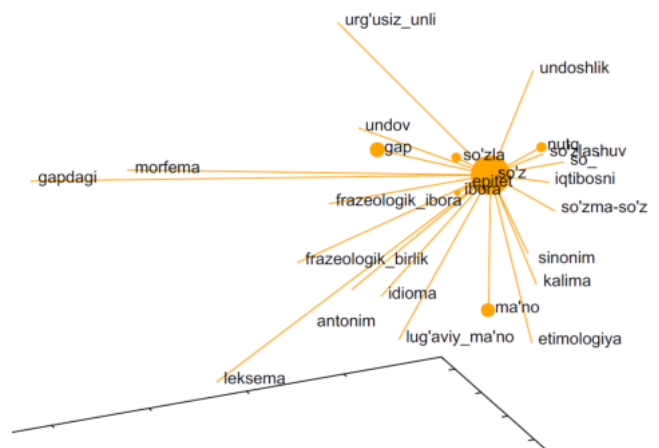


Fig. 6.   Three dimensional visualization of word embedding vectors, in the vicinity of the lemma *so'z*.

---

[9] https://www.juls.savba.sk/semä/?lang=uz

There are three different models available: the model trained on the combination of part of speech and lemmas, the model trained on word forms, and a FastText model The



Fig. 8. List of nearest vectors to the result of an arithmetic expression qirol – er + ayol.

models are also provided as downloadable datasets in text GenSim format.

The links to external resources are (together with the hyperlink symbols): link to the Aranem Uzbecicum Minus corpus; link to the Google search for the word, restricted to the *.uz* top level domain; link to English language Wiktionary entry for the word. This allows quick verification of the usage of the word, mostly for lexicographic purposes.



Fig. 7. List of nearest vectors to the lemma *so'z*. The columns contain sine similarity, lemma, absolute frequency in the corpus, link to the *Araneum Uzbecicum Minus* corpus, link to Google Search, link to Wiktionary.

The interface supports two other modes of operation: if there are two words (separated by a space) in the input, it displays the sine similarity between them.

The other mode gives access to a basic vector arithmetics, supporting addition and subtraction. This allows easy evaluation of the well known vector transfer (the famous king – man + woman = queen example, see [14] with many applications in linguistic and sociolinguistic research. In Fig. 8. We provide an Uzbek equivalent, the nearest vectors to the expression qirol - er + ayol correspond to the words malika and qirolicha, meeting our expectations.

## XI. Conclusion and Further Work

Based on the feedback of the early users, we would like to perform further processing in the foreseeable future as follows:

• As the corpus contains some amount of non-Uzbek (Mostly English and Russian) text fragments, a sentence-level language filtration using the methodology suggested by [28] will be applied.

• Sketch and term grammars for the Sketch Engine will be crated.

• Based on the Apertium morphological lexicon, we'll try to train rules for the CSTlemma lemmatizer, so that lemmas for the OOV lexical items could be guessed.

Besides the availability of our new corpus for on-line querying at our Corpus Portal, we can also provide the processed corpus source data (for non-commercial purposes) under an open license. We hope this this might improve the situation with availability of large-scale language resources for the Uzbek language.

The full Gigaword version of the new corpus will be published on-line in a foreseeable future.

### References

[1] N. Abdurakhmonova, U.Tuliyev, A. Gatiatullin. Linguistic functionality of Uzbek Electron Corpus: uzbekcorpus.uz // 2021 International Conference on Information Science and Communications Technologies (ICISCT), Tashkent, Uzbekistan, pp. 1-4, (2021).

[2] S. S. Avezov Development of a Parallel Corpus of the Uzbek and Russian Languages // Vital Annex: International Journal of Novel Research in Advanced Sciences (IJNRAS) Volume: 01, Issue: 05, (2022).

[3] V.Baisa & V.Suchomel. Turkic language support in Sketch Engine. // 3rd International Conference on Computer Processing in Turkic Languages (TURKLANG 2015).

[4] A.Barbaresi. Challenges in web corpus construction for low-resource languages in a post-BootCaT world. 6th Language & Technology Conference, Less Resourced Languages special track, Poznan, Poland. pp. 69-73 (2013).

[5] Baroni, M., & Bernardini, S. BootCaT: Bootstrapping Corpora and Terms from the Web. In Proceedings of LREC 2004, pp. 1313-1316 (2004).

[6] V. Benko Aranea: Yet Another Family of (Comparable) Web Corpora. In Petr Sojka, Aleš Horák, Ivan Kopeček and Karel Pala (Eds.): Text, Speech and Dialogue. 17th In-ternational Conference, TSD 2014,

Brno, Czech Republic, September 8-12, 2014. Pro-ceedings. LNCS 8655. Springer International Publishing Switzerland (2014)

[7]  G. Radovan: Word Embedding Based on Large-Scale Web Corpora as a Powerful Lexicographic Tool. In: Rasprave: Časopis Instituta za hrvatski jezik i jezikoslovlje 46, br. 2 (2020): 603-618

[8]  B.Jongejan and H.Dalianis. Automatic training of lemmatization rules that handle morphological changes in pre-, in- and suffixes alike. In Proceedings of the Joint Con-ference of the 47th Annual Meeting of the ACL and the 4th International Joint Confer-ence on Natural Language Processing of the AFNLP. Suntec, Singapore: Association for Computational Linguistics, pp. 145-153 (2009).

[9]  T.Khanna, J. N. Washington, F. M.Tyers, S.Bayatlı, D. G.Swanson, T. A.Pirinen, I.Tang, and Alòs i Font, H. Recent advances in Apertium, a free/open-source rule-based machine translation platform for low-resource languages. Machine Translation (2021)

[10]  A.Kilgarriff, V.Baisa, J. Bušta, M.Jakubíček, V.Kovář, J.Michelfeit, P.Rychlý, and V.Suchomel. The Sketch Engine: Ten Years On. Lexicography 1(1):7–36 DOI: 10.1007/s40607-014-0009-9 (2014).

[11]  N.Ljubešić and F. Klubička. {bs,hr,sr}WaC – Web corpora of Bosnian, Croatian and Serbian. Proceedings of the 9th Web as Corpus Workshop (WaC-9). Gothenburg, Sweden (2014).

[12]  R.McDonald, J.Nivre, Y.Quirmbach-Brundage, Z.Goldberg, D.Das, K.Ganchev, K.Hall, S.Petrov, H. Zhang, 0.Täckstrom, C.Bedini, Bertomeu Castelló, N., Lee, J. Universal Dependency Annotation for Multilingual Parsing. In Proceedings of ACL (2013).

[13]  J.Michelfeit, J.Pomikálek, V.Suchomel. Text Tokenisation Using unitok. In 8th Workshop on Recent Advances in Slavonic Natural Language Processing, Brno, Tribun EU, pp. 71-75 (2014).

[14]  T.Mikolov, K ai Chen, G.Corrado, J. Dean, 2013. Efficient Estimation of Word Representations in Vector Space. Proceedings of Workshop at ICLR. Université de Montreal. Scottsdale.

[15]  T.Mikolov, E.Grave, P.Bojanowski, Ch.Puhrsch, A.Joulin. 2018. Advances in Pre-Training Distributed Word Representations. Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018). European Language Resources Association. Miyazaki.

[16]  M.Musaev, S.Mussakhojayeva, I.Xujayorov, Y.Khassanov, M.Ochilov, A.Varol. USC: An Open-Source Uzbek Speech Corpus and Initial Speech Recognition Experiments. (2021)

[17]  P.Paikens. Deep Neural Learning Approaches for Latvian Morphological Tagging. In Human Langiage Technologies – The Balitc Perspective (2016).

[18]  J. M. Patel. Introduction to Common Crawl Datasets. In Getting Structured Data from the Internet, pp 277–324. Apress (2020)

[19]  J.Pomikálek. Removing boilerplate and duplicate content from web corpora. PhD the-sis, Masaryk University, Faculty of informatics, Brno, Czech Republic (2011).

[20]  M. O. Rabin (1981). Fingerprinting by Random Polynomials. Center for Research in Computing Technology, Harvard University. Tech Report TR-CSE-03-01 (1981).

[21]  R.Řehůřek, P.Sojka. Software framework for topic modelling with large corpora. 2010. Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks.

[22]  P.Rychlý. Manatee/Bonito – A Modular Corpus Manager. In 1st Workshop on Recent Advances in Slavonic Natural Language Processing. Brno: Masaryk University, 2007. p. 65–70 (2007).

[23]  R.Schäfer and Bildhauer. Building Large Corpora from the Web Using a New Efficient Tool Chain. In Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12). 486–493 (2012).

[24]  H.Schmid. Probabilistic Part-of-Speech Tagging Using Decision Trees. Proceedings of International Conference on New Methods in Language Processing, Manchester (1994).

[25]  D.Spoustová. J.Hajič, J.Raab, M. Spousta. Semi-Supervised Training for the Averaged Perceptron POS Tagger. In Proceedings of the 12th Conference of the Eu-ropean Chapter of the ACL (EACL), pages 763-

771, Athens, Greece, March. Associa-tion for Computational Linguistics (2019).

[26]  M.Straka, J.Hajič, J.Straková. UDPipe: Trainable Pipeline for Processing CoNLL-U Files Performing Tokenization, Morphological Analysis, POS Tagging and Parsing. In Proceedings of the Tenth International Conference on Language Resources and Evalua-tion (LREC), Portorož, Slovenia (2016)

[27]  V.Suchomel, J.Pomikálek, Efficient Web Crawling for Large Text Corpora. In Adam Kilgarriff, Serge Sharoff. Proceedings of the seventh Web as Corpus Workshop (WAC7). Lyon, pp. 39–43 (2012).

[28]  V. Suchomel. Discriminating Between Similar Languages Using Large Web Corpora in the Proceedings of Recent Advances in Slavonic Natural Language Processing, RASLAN (2019).

[29]  M.Sharipov, J.Mattiev, J.Sobirov, R.Baltayev. Creating a Morphological and Syntactic Tagged Corpus for the Uzbek Language // https://ceur-ws.org/Vol-3315/paper10.pdf (2022)

[30]  A.S. Sobirovich (2022). Development of a parallel corpus of the uzbek and russian languages. Vital Annex: International Journal of Novel Research in Advanced Sciences, 1(5), 152–155.

[31]  M. Tursunov (2023). Software of the national corpus of the uzbek language. International Journal of Advance Scientific Research, 3(10), 190-199.

[32]  M.S. Tursunov (2022). Development of software for a text corpus in uzbek // Native Languages and Cultures in the Modern Changing World. №1, 62-70 p.

[33]  M.Volk (2002). Using the web as corpus for linguistic research. 10.5167/uzh-20339.

[34]  G. Eason, B. Noble, and I. N. Sneddon, "On certain integrals of Lipschitz-Hankel type involving products of Bessel functions," Phil. Trans. Roy. Soc. London, vol. A247, pp. 529–551, April 1955. (references)

[35]  J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.

[36]  I. S. Jacobs and C. P. Bean, "Fine particles, thin films and exchange anisotropy," in Magnetism, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271–350.

[37]  K. Elissa, "Title of paper if known," unpublished.

[38]  R. Nicole, "Title of paper with only first word capitalized," J. Name Stand. Abbrev., in press.

[39]  Y. Yorozu, M. Hirano, K. Oka, and Y. Tagawa, "Electron spectroscopy studies on magneto-optical media and plastic substrate interface," IEEE Transl. J. Magn. Japan, vol. 2, pp. 740–741, August 1987 [Digests 9th Annual Conf. Magnetics Japan, p. 301, 1982].

[40]  M. Young, The Technical Writer's Handbook. Mill Valley, CA: University Science, 1989.

[41]  B. R. Mengliyev, S.Hamroyeva, O.Abdullayeva (2023). Scopus-based bibliometric analysis on corpus linguistics for the period of 2017-2021. In E3S Web of Conferences (Vol. 413, p. 03008). EDP Sciences.

[42]  B.Mengliyev, S.Shahabitdinova, S.Khamroeva, S.Gulyamova, A.Botirova. (2021). The morphological analysis and synthesis of word forms in the linguistic analyzer. Journal of Language and Linguistic Studies, 17(1), 558-564.

[43]  K. R.Abdurasulovich, M. B. Rajabovich. (2019). The Role of the Parallel Corpus in Linguistics, the Importance and the Possibilities of Interpretation. International Journal of Engineering and Advanced Technology, 8(5), 388-391.

[44]  E.B.Boltayevich, S.S.Samariddinovich, K.S.Mirdjonovna, E.Adali, X.Z. Yuldashevna. Pos Taging of Uzbek Text Using Hidden Markov Model. UBMK 2023 - Proceedings: 8th International Conference on Computer Science and Engineering, 2023, P. 63–68.

[45]  E.B.Boltayevich, E.Adali, K.S.Mirdjonovna, X.Z.Yuldashevna, N.Uktamboy o'g'li, X.The Problem of Pos Tagging and Stemming for Agglutinative Languages (Turkish, Uyghur, Uzbek Languages). UBMK 2023 - Proceedings: 8th International Conference on Computer Science and Engineering, 2023, P. 57–62.