



Automatic Language Identification and Relatedness Mapping

Sonia Cromp

Goals

- Language identification
- Language relatedness
 - Ancestral
 - Geographic
 - Cultural
 - Common writing system

عین و لست قویه آه تو میستادی شامونیه ماداییه
دآ اویباقه سایم کاده، قلم فنه این تیشان دن شاه
دایسیده، دآ پیغمه سلیم ما قدّتہ و بک
اونچی مونیه کاف سلپسنه سیچه مو ریگه
یاسام قوقیه پیسیه منی یورکدا آئی میبیه.



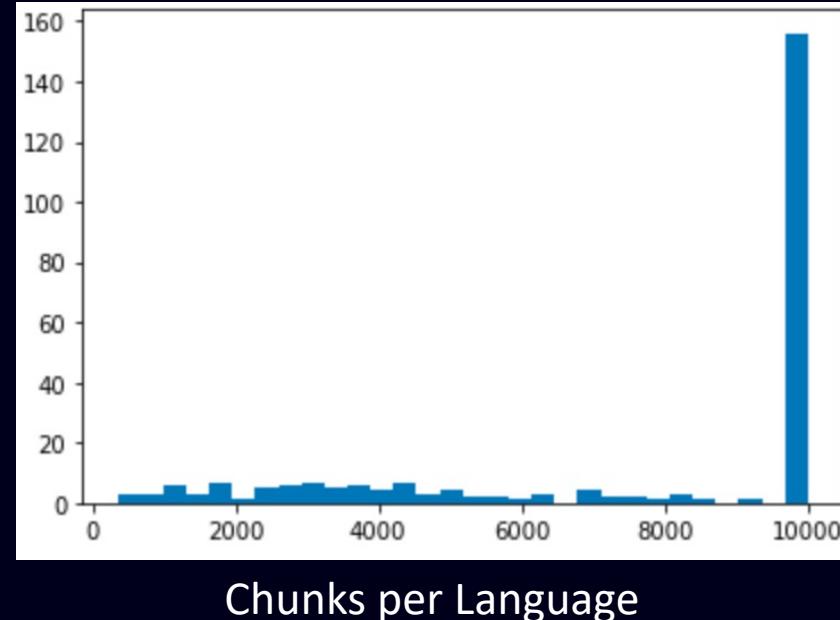
Nacionalni sabor
socijalističkih
snaga
BIROKRATSKI PUC
U BARANJI
O NACIONALNOJ
STRUKTURI U

VLADO GOTOVAC
**TRAŽIMO
DOKAZE**

За овој врсник наје замје постава пре-
тврда, која је сјајнији почетак привреде. Баде ре
на нападе индустрије биоса и најбоље изјаве бара-
не су уједињене у једном месту.

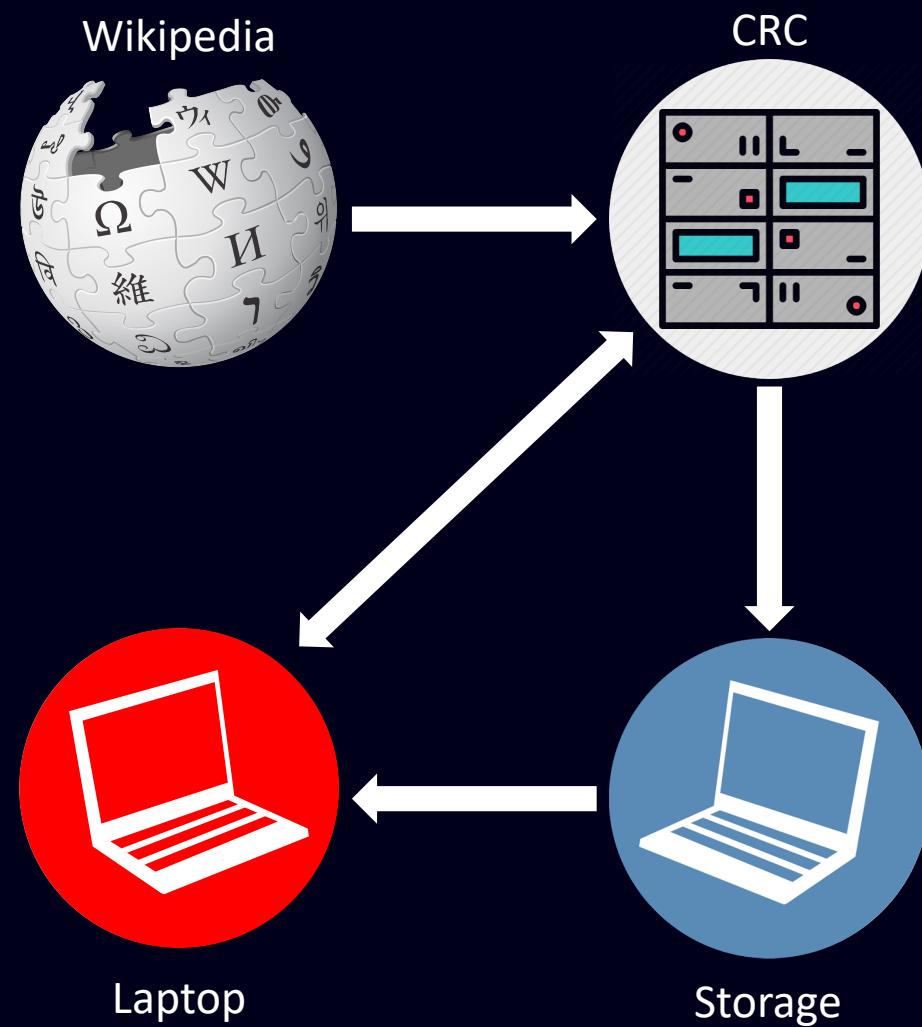
Dataset

- 248 languages' Wikipedia article texts
- Cleaning: Remove punctuation, non-text characters
- Result: Random 500-character “chunks”



Data gathering

1. Networking
2. Download server dumps
3. Remove formatting, clean text
4. Chunking and shuffling



Data gathering

Before (random section of Qaraqalpaqsha)

'<timestamp>2020-05-01T08:16:04Z</timestamp>\n', '<contributor>\n', '<ip>185.163.26.25</ip>\n', '</contributor>\n', '<model>wikitext</model>\n', '<format>text/x-wiki</format>\n', '<text bytes="7440" xml:space="preserve">\"\'A\'jiniyaz Qosıbay ulı\"\' (a\'debiy laqabı \"'Ziywar\"') - XIX a\'sirdegi qaraqalpaq klassik a\'debiyatının\' en\' o\'rnekli wa\'killerinin\' biri. Ol do\'retiwshiliginde o\'zine ta\'n o\'zgeshelikke iye, og\'ada talantlı, oqmışlı, medreseni ayrıqsha tamamlag\'an aqun, ulama, ko\'rkem so\'z sheberi sıpatında basqlardan ayrılp turatug\'in uqıplılıq penen o\'z da\'wirinin\' progressiv ...',

After (one English chunk)

'alism among other libertarian socialist economic theories As anarchism does not offer a fixed body of doctrine from a single particular worldview many anarchist types and traditions exist and varieties of anarchy diverge widely One reaction against sectarianism within the anarchist milieu was anarchism without adjectives a call for toleration and unity among anarchists first adopted by Fernando Tarrida del Mármol in response to the bitter debates of anarchist theory at the time Belief in poli\n'

248 languages included

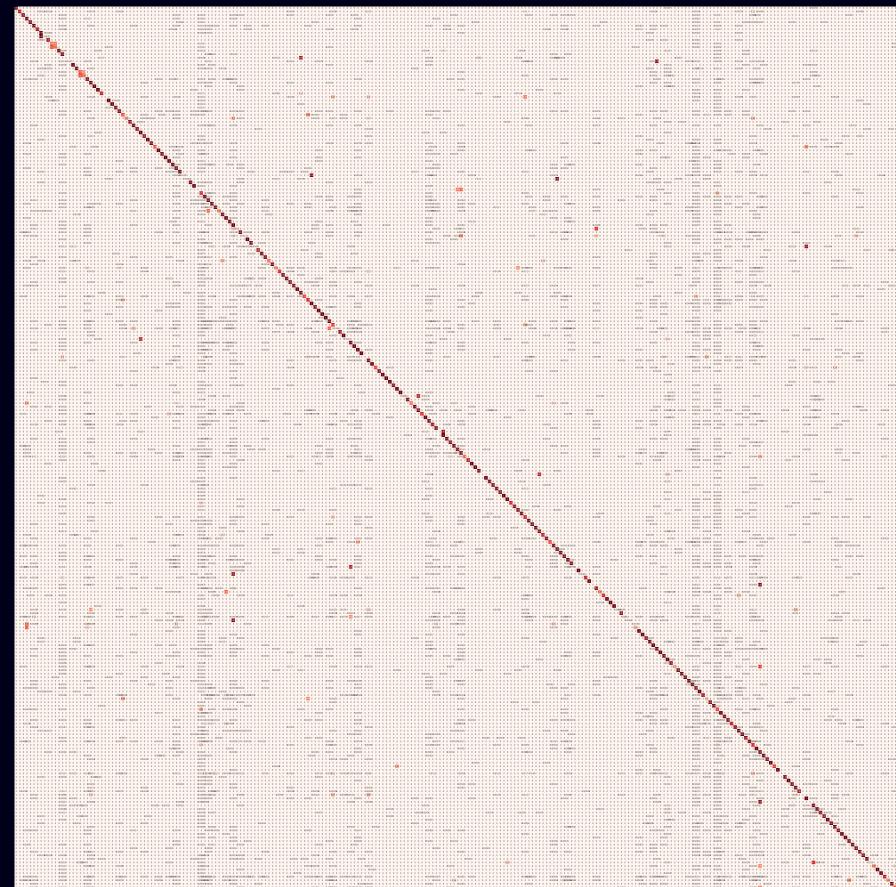
Kreyòl ayisyen, Limburgs, Igbo, tetun, 吴语, коми, الادارجة latviešu, la .lojban., મોજપુરી, мокшень, Zeêuws, o'zbekcha/ўзбекча, تۆركچى, norsk nynorsk, português, Bahasa Melayu, અવધી, 中文, тыва дыл, svenska, башкортса, ଭାଲ୍ଯୁଗ୍ଭୁ, Tagalog, arpetan, slovenščina, Afrikaans, Jawa, Yorùbá, Mirandés, જુજરાતી, Frysk, Māori, Xitsonga, sicilianu, italiano, Aymar aru, Gagauz, କୁର୍ରା, davvisámegiella, føroyskt, ତେଉଳୁ, brezhoneg, català, Адыгэбзэ, dansk, 贛語, Pälzisch, Nāhuatl, Novial, qırımtatarca, Simple English, dolnoserbski, ବାହୁଦ୍ରାଷ୍ଟ୍, Oromoo, અસમীয়া, English, lumbaart, বাংলা, , hrvatski, Nouormand, ନୀତି, ଓଡ଼ିଆ, русский, 한국어, interlingua, Papiamentu, қарабачай-малкъар, Кыргызча, Kotava, български, azərbaycanca, vepsän kel', Alemannisch, Gaelg, Plattendütsch, livvinkarjala, Soomaaliga, Piemontèis, Sunda, Ænglisc, galego, Hausa, Mìng-děng-ngū, יִדִיש, emiliàn e rumagnòl, پنجابی, авар, walon, română, Picard, kurdî, مصرى, سُونِي, suomi, මුණ්ඩාවාවා, kaszëbsczi, Bahasa Indonesia, ବୈଜ୍ଞାନିକ, vèneto, eesti, sardu, Kapampangan, удмурт, олык марий, 日本語, Ελληνικά, саха тыла, Runa Simi, Deitsch, bosanski, ନାଥାନ୍ତାନ୍ତା, español, سندھی, қырык мары, aragonés, Ilokano, Banjar, Kinyarwanda, 客家語/Hak-kâ-ngî, српски / srpski, Wolof, Türkçe, беларуская, Scots, Fiji Hindi, Чӑвашла, مازرۇنى, anarâškielâ, хальмг, hujjebek, West-Vlams, Галгай, ဘာသာ မန်, Zazaki, occitan, ମାର୍ଗାଲ୍ଲୁରୀ, Avañe'ẽ, čeština, Tok Pisin, Sesotho sa Leboa, Kiswahili, Boarisch, ქართული, Seeltersk, Latina, Ido, ମୈଥୁଣୀ, मराठी, norsk, эрзянь, srpskohrvatski / српскохрватски, Lingua Franca Nova, Аԥସା, ତୁଳୁ, Taqbaylit, Winaray, Ripoarisch, rumantsch, мæथିଲୀ, русиньскый, kernowek, Chi-Chewa, лакку, Sakizaya, magyar, нохчийн, اردو, Volapük, Deutsch, Ligure, संस्कृतम्, лезги, asturianu, தமிழ், नेपाली, euskara, اسرائیلی Interlingue, Ladino, Napulitano, ଡୋଟେଲୀ, қазақша, Malagasy, Gàidhlig, lietuvių, Cymraeg, Qaraqalpaqsha, हिन्दी, перем коми, فارسی, Ślůnski, Türkmençe, slovenčina, Nederlands, isiZulu, Lëtzebuergesch, corsu, татарча/tatarça, Cebuano, íslenska, монгол, Basa Bali, Acèh, ବିଶ୍ୱାସ୍ମିଯା ମନିପୁନ୍ନୀ, Minangkabau, ئۇيغۇرچە / Uyghurche, Luganda, chiShona, estremeñu, furlan, Setswana, پښتو, Malti, Nordfriisk, עברית Kabiyé, hornjoserbsce, shqip, Diné bizaad, Pangasinan, polski, македонски, Ќумоџија, Tiếng Việt, Արեւմտահայերէն, ठन್ನಡ, ଗୋଁଚୀ କାକଣୀ / Gõychi Konknni, Ирон, lingála, тоҷикӣ, ҳӯҷӣ, Bahasa Hulontalo, українська, Gaeilge, नेपाल भाषा, ສັດໝີ, français, Bikol Central, isiXhosa, Esperanto, العربية

Language ID

Language Identification

- Count Vectorizer on character bigrams
- Naïve Bayes
- 90.6% accuracy

True Language



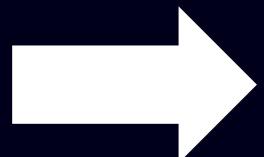
Predicted Language

Relatedness mapping

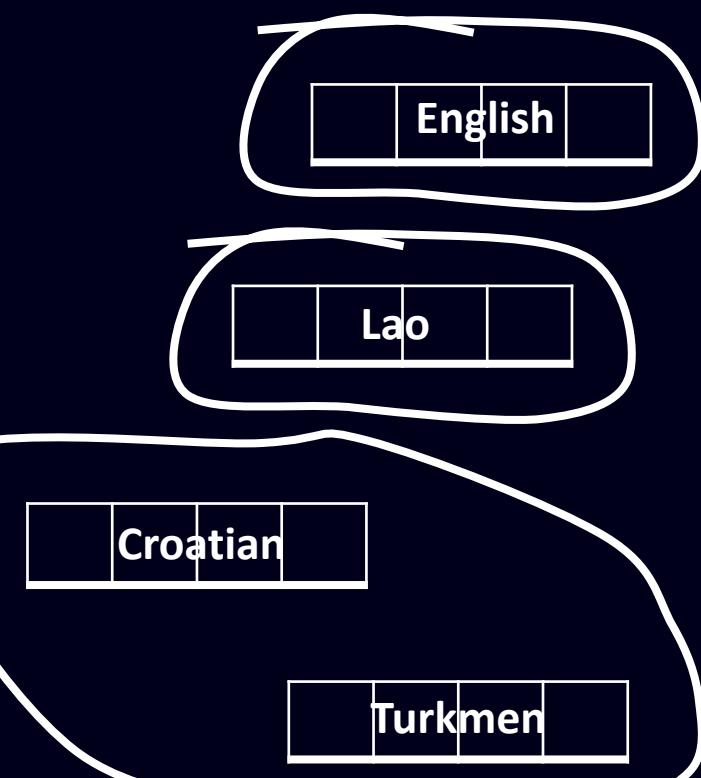
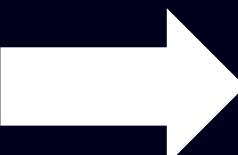
Mapping Relatedness

- k -means clustering, $k=150$

True Language	English	Croatian	Lao	Turkmen
English	8	1	0	1
Croatian	1	6	0	3
Lao	1	0	9	0
Turkmen	0	3	1	6



8	1	0	1
1	6	0	3
1	0	9	0
0	3	1	6



Clusters

- Swahili, Tsonga, Xhosa, Chewa, Kinyarwanda (Bantu)
- French, Norman, Picard (Oil)
- Persian, Gilaki, Mazanderani (Western Iranian)
- Central Bikol, Cebuano, Javanese, Pangasinan, Tagalog (Malayo-Polynesian),
Tok Pisin (English creole)
- Alemannic, Ripuarian, Pennsylvania German, Palatine German (High German)
- Romanian (Balkan Romance), Silesian (West Slavic), Kotava, Lojban (constructed)

Anonymizing writing systems

- For all chunked text in the language, rank all characters (except space) by commonality
- Replace space with 0, most common character with 1, second most with 2, etc up to 256
- Count Vectorizer on character bigrams
- Naïve Bayes
- 99% accuracy

con esta renas la imperos brites

```
\x0c\x06\x05\x00\x02\x03\t\x01\x00\x08\x02\x05\x01\x03\x00\x07\x01\x00\x04
\x0e\r\x02\x08\x06\x03\x00\x12\x08\x04\t\x02\x03
```

Conclusions and future work

- Geographical > Ancestral (eg Romanian, Tok Pisin)
- Writing systems aren't helpful, or rarer characters aren't helpful

Future work

- Understand anonymized writing system results
- Explore hyperparameters