

**STATISTICAL APPROACH FOR
RECOGNITION AND
DETECTION OF PHONETIC
BORROWINGS IN CHINESE
(CASE STUDY OF THE
RUSSIAN PROPER NAMES)**

Kirill Semenov

HSE – Moscow

eReL: May 11th, 2019

OBJECT: IS IT
IMPORTANT?..

There are other more popular spheres in modern Chinese lexicography:

- Semantic borrowings: Superman > 超人 (chāorén = “exceed”+ “man”)
- Loanword blending: ballet > 芭蕾舞 (bālěiwǔ = phonetic loan + “dance”)
- “Lettered words”: 卡拉OK (kǎlā ou kei – “karaoke”);
三G手机 (sān Gē shǒujī – “three” + “Generation [English]” + “mobile phone”)
- Modern Internet neologisms
 - 同志 (tóngzhì) – “comrade” > “homosexual”
 - 1314 (yī sān yī sì) ≈ yī sheng yī shì = “forever”

OBJECT: IT IS
IMPORTANT!

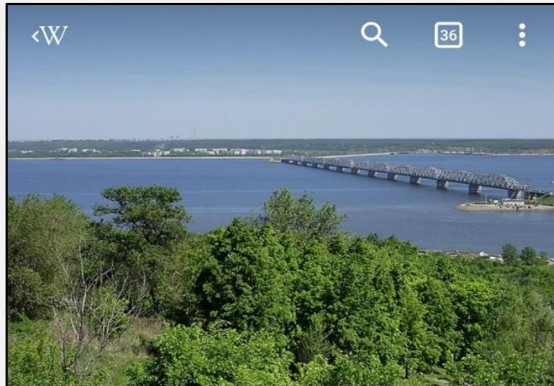
For the fundamental research:

- phonetic adaptation in a language of different phonetic inventory

For the natural language processing (NLP) purposes:

- PoS-tagger
- NER (Named Entity Recognition)
- MT (Machine translation)
- etc

OBJECT



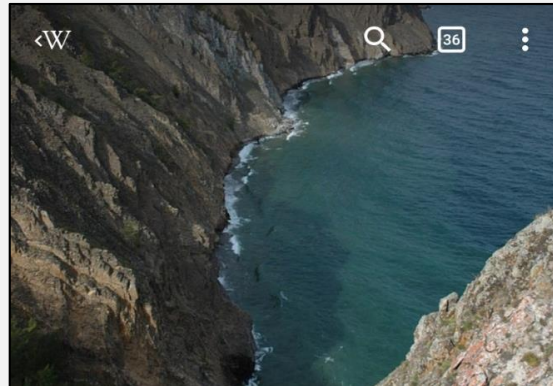
伏尔加河

— 添加标题描述

伏尔加河（**俄语**：Волга，又译窝瓦河），位于**俄罗斯**西南部，全长3,692公里，是**欧洲**最长的河流^[2]，也是世界最长的**内流河**^[2]，流入**里海**。

快速预览: 国家, 城市 ...

伏尔加河是欧洲**流域**最广以及流量最大的河流，流经**欧洲俄罗斯**，是代表型的**俄罗斯河**。



贝加尔湖

俄罗斯湖泊

贝加尔湖（**俄语**：озеро Байка́л，**罗马化**：Ozero Baykal，IPA：[ɒzʲɪrə bɐjˈkal]；**布里亚特语**：Байгал нуур，拉丁转写：Baygal nuur，**蒙古语**：Байгал нуур；意思是“自然之湖”^[3]；一说名称来源于“贝音嘎 嘎拉”（蒙古语意为不灭的火焰）^[4]）。**汉朝**人称之为“翰海”^[5]，五胡十六国时北朝叫“于已尼大水”，隋唐叫“小海”，18世纪初期的《**异域录**》称之为“**柏海儿湖**”，



莫斯科

俄罗斯首都

莫斯科（**俄语**：Москва，**罗马化**：Moskva，IPA：[mesˈkva]）是**俄罗斯首都**与最大都市、以及**莫斯科州**首府，为**俄罗斯**全国政治、经济、科学、文化及交通的中心。面积2,510平方公里，与**莫斯科州**和**卡卢加州**接壤。**城区**人口约1200万，是**欧洲**人口第二多的城市，仅次于**伊斯坦布尔**，占俄罗斯总人口的1/10。



列夫·托尔斯泰

俄国作家

列夫·尼古拉耶维奇·托尔斯泰（**俄语**：Лев Николаевич Толстой，**拉丁化**：Lev Nikolayevich Tolstoy；1828年9月9日（**儒略历**8月28日）—1910年11月20日（**儒略历**11月7日）），**俄国**小说家、哲学家、政治思想家，也是**非暴力**的**基督教无政府主义者**和教育改革家。他是在**托尔斯泰**这个**贵族**家族中最有影响力的一位。

PROBLEMS: PHONETIC ASPECT

A big difference in phonetic inventories of the SAE languages and the Chinese

SAE	Chinese
Voiced VS Voiceless consonants	Aspirated VS Unaspirated consonants
Stress	Tone

A big level of homonymy in Chinese

- shí: 时 (a while), 十 (ten), 石 (stone)...
- shì: 是 (to be), 市 (town), 事 (case), 试 (to try), 世 (generation), 示 (to demonstrate)...

PROBLEMS: GRAPHICAL ASPECT

Absence of spaces between words

- 本宪章不得认为授权联合国干涉在本质上属于任何国家国内管辖之事件，且并不要求会员国将该项事件依本宪章提请解决；但此项原则不妨碍第七章内执行办法之适用。

No set of symbols used specifically for the phonetic loanwords

- 马里 (mǎlǐ) = “Horse mile/Inside horse” > Mali (African country) / Mary (city in Turkmenistan)

Ambiguity of the characters

- 乐 (yuè) = “music” / (lè) “happy”
- 了 (le) = PERF / (liǎo) = “to understand”

Our Hypothesis:

There is a particular pattern of transliteration for the Russian words in Chinese.

Our Aim:

To find it!

OUR TASKS

1. To check whether the transliteration trends of English and German are applicable to the Russian loanwords
 - Case study: to check whether there is influence of a Chinese-Russian pidgin on the current Russian loanwords' adaptation
2. To compare the prescribed transliteration rules (Xinhua) to the real data (Wikidata, dictionaries)
3. To analyze the cases of the partial semantic translation of the Named Entities (Wikidata, dictionaries)
4. To compare the most frequent N-grams in the Russian NE, Chinese NE and the reference corpus (Wikidata)

I. PHONETIC ADAPTATION OF THE RUSSIAN WORDS IN CHINESE: THE OT APPROACH

PHONETIC ADAPTATION OF THE RUSSIAN WORDS IN CHINESE: THE OT APPROACH

Based on the study of the English, German and Italian loanword adaptation in Chinese (Miao 2005)

- Alternation of consonant phonemes
- Transformation of consonant clusters

Made in the Optimality Theory paradigm

PHONETIC ADAPTATION
OF THE RUSSIAN WORDS
IN CHINESE: RESULTS

The main principles are confirmed on the Russian data:

- The crucial consonant feature is MANNER:
 - $s \Rightarrow s$ (75%) $\gg \zeta$ (20%) $\gg \mathfrak{s}$ (5%);
 $*d\mathfrak{z}$
 - $t \Rightarrow t^h$ (70%) $\gg t$ (20%) $\gg *! tsj$
(10%); $*n$
- The C deletion rate (Coda position):
 - Liquids \gg plosives \gg fricatives \gg
/m, n/

PHONETIC ADAPTATION OF THE RUSSIAN WORDS IN CHINESE: RESULTS

Hard VS Soft feature of Russian consonants

In general – does not affect the transformation

Soft velars – c, ʃ, ç:

- In Russian – palatal place and a fricative-like release
- In Chinese – substituted by affricates tʃʲ tsʲ ʈʃʲ
- detected in English as well, but more consistent in Russian

Soft dentals - dʲ , tʲ:

- In Russian - almost affricates tʃʲ tsʲ
- In Chinese - substituted by the corresponding affricates

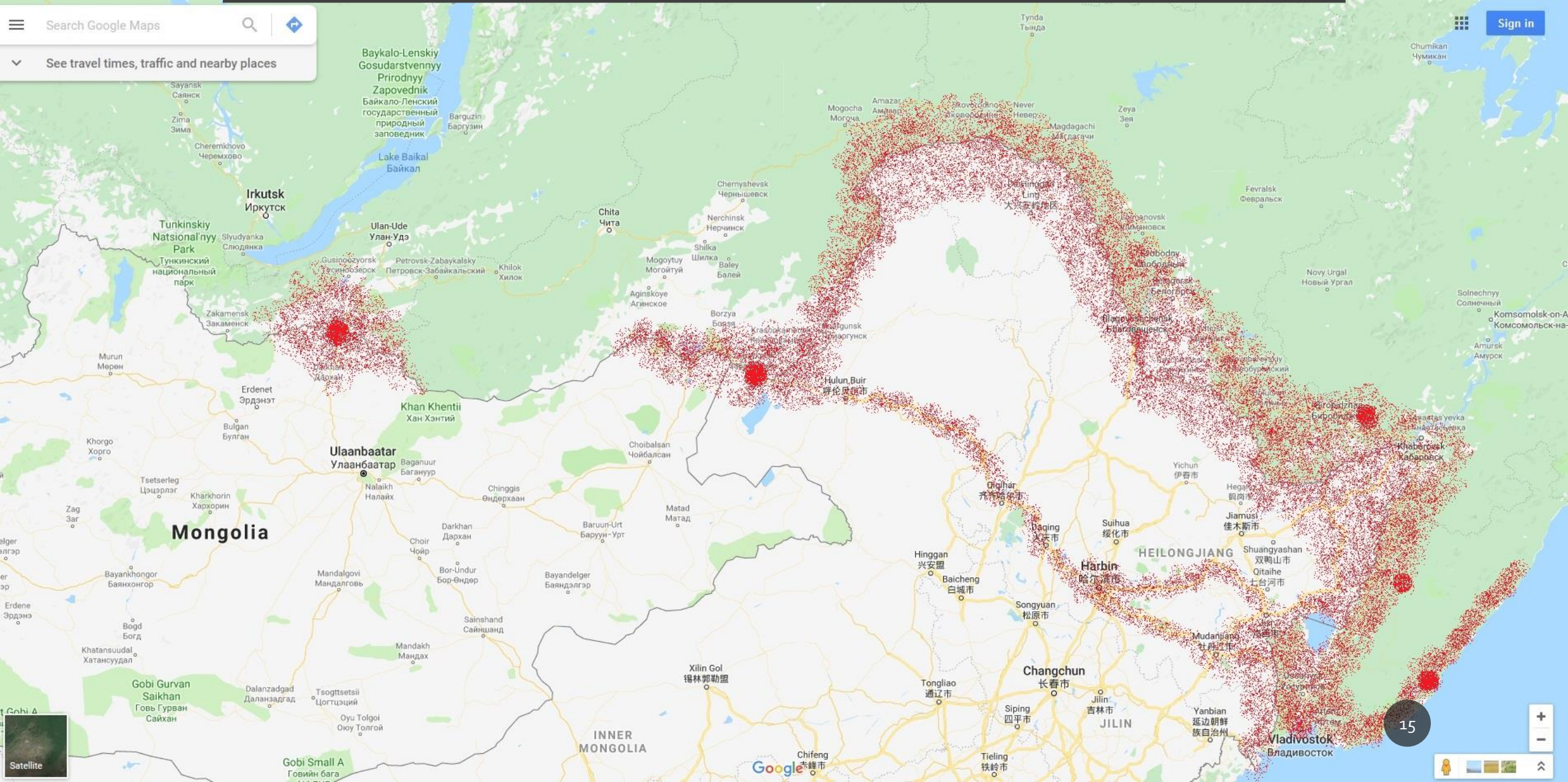
CASE STUDY: CHINESE-RUSSIAN PIDGIN

CASE STUDY: CHINESE-RUSSIAN PIDGIN

- Russian has a different history of contact with Chinese
- Pidgin has existed since the end of XVIII cent.



CASE STUDY: CHINESE-RUSSIAN PIDGIN



CASE STUDY: HYPOTHESIS

The Chinese-Russian pidgin can influence:

- The phonetic appearance of the exact “old” loanwords
- The phonetic adaptation strategy for the new Russian loanwords

We can state it if:

- The phonetic adaptation differs from the usual Chinese one
- The phonetic adaptation is made according to the dialect phonology

CASE STUDY: RESULTS

The majority of the dialectal features from the standard Chinese is not observed in the loanwords

If it is, it can be usually explained by the OT assumptions:

- 哈拉嗦 hālāsuó instead of halaşuo – predicted by the variation within one MANNER type

CASE STUDY: RESULTS

Two aspects likely to be transferred from pidgin:

- “dza-adaptation” of nouns – a very “pidginish” feature

Russian	Pidgin	Northern Dialects	Modern Chinese
Купец	kupe-dza	谷瘰-子 gǔbiě-zī*	谷瘰-子 gǔbiě-zī*
Халат	Hala-dza	哈拉-子 hālā-zī*	哈拉-呢 hālā-nǐ

- The adaptation of a sound [z] is an affricate [dʒ] or [dʒ̥], which differs from the predicted by OT adaptation (same for pidgin and for contemporary Russian loans)

	Russian	Pidgin	Modern Chinese
XX cent. Borrowing	КОЛХОЗ	kaxódʒə	科尔火支 kē'ěrhǔōzhī*
New Borrowing (NE)	Рязань	---	梁赞 liángzàn*

*by official “pinyin” transcription, “z” = [dʒ], “zh” = [dʒ̥]

2. PRESCRIPTIONS AND USE

XINHUA ALGORITHM

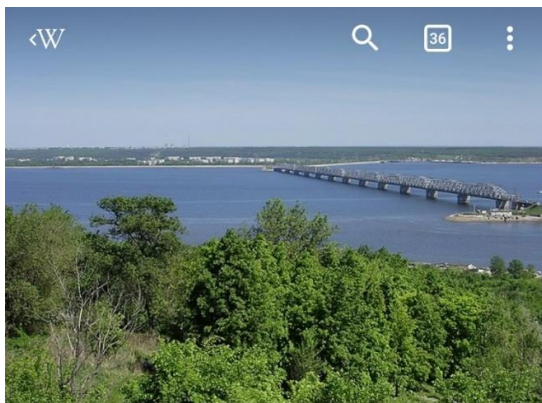
Python code

- Regular expressions
- “if”- “else” constructions
- Pandas library
 - .csv table with prescriptions

ALGORITHM APPLICATION

Wikidata

DATA: WIKIPEDIA OBJECTS



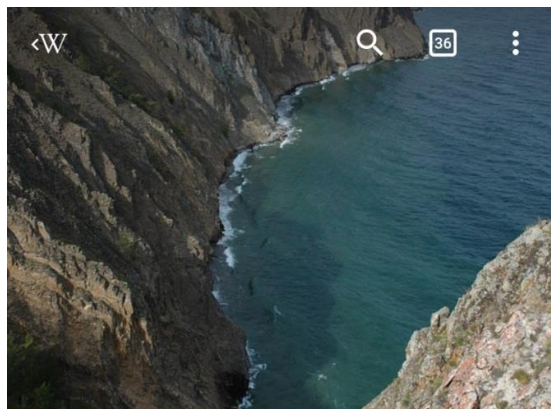
伏尔加河

— 添加标题描述

伏尔加河（**俄语**：Волга，又译窝瓦河），位于**俄罗斯**西南部，全长3,692公里，是**欧洲**最长的河流^[2]，也是世界最长的**内流河**^[2]，流入里海。

快速预览: 国家, 城市 ...

伏尔加河是欧洲**流域**最广以及流量最大的河流，流经**欧洲俄罗斯**，是代表典型的**俄罗斯河**



贝加尔湖

俄罗斯湖泊

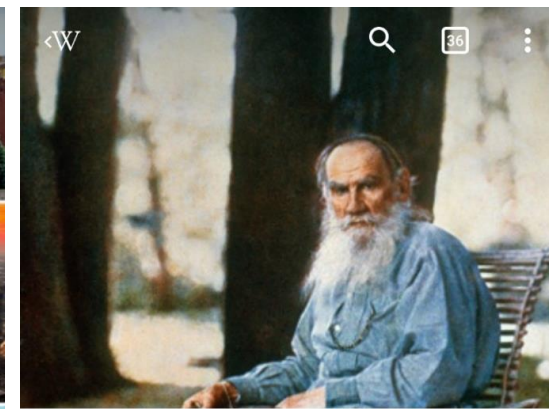
贝加尔湖（**俄语**：озеро Байка́л，**罗马化**：Ozero Baykal，IPA：[ˈozɪrə beɪˈkal]；**布里亚特语**：Байгал нуур，拉丁转写：Baygal nuur，**蒙古语**：Байгал нуур；意思是“自然之湖”^[3]；一说名称来源于“贝音嘎 嘎拉”（蒙古语意为不灭的火焰）^[4]）。**汉朝**人称之为“翰海”^[5]，五胡十六国时北朝叫“于已尼大水”，隋唐叫“小海”，18世纪初期的《**异域录**》称之为“柏海儿湖”，



莫斯科

俄罗斯首都

莫斯科（**俄语**：Москва，**罗马化**：Moskva，IPA：[mesˈkva]）是**俄罗斯首都**与最大都市、以及**莫斯科州**首府，为**俄罗斯**全国政治、经济、科学、文化及交通的中心。面积2,510平方公里，与**莫斯科州**和**卡卢加州**接壤。**城区**人口约1200万，是**欧洲**人口第二多的城市，仅次于**伊斯坦布尔**，占俄罗斯总人口的1/10。



列夫·托尔斯泰

俄国作家

列夫·尼古拉耶维奇·托尔斯泰（**俄语**：Лев Николаевич Толстой，**拉丁化**：Lev Nikolayevich Tolstoy；1828年9月9日（**儒略历** 8月28日）—1910年11月20日（**儒略历** 11月7日）），**俄国**小说家、哲学家、政治思想家，也是**非暴力**的**基督教无政府主义者**和教育改革家。他是在**托尔斯泰**这个**贵族**家族中最有影响力的一位。



Main Page

Discussion

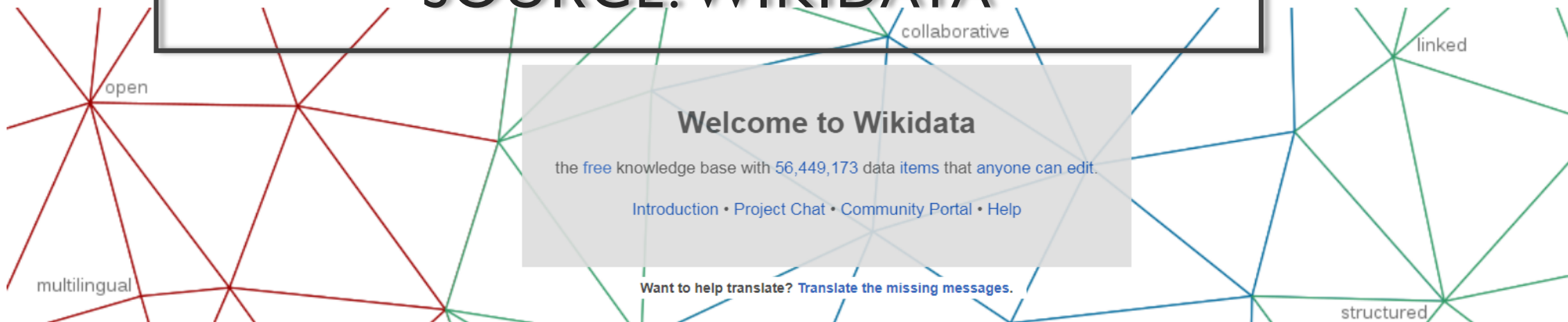
SOURCE: WIKIDATA

Read

View source

View history

Search Wikidata



Welcome to Wikidata

the free knowledge base with 56,449,173 data items that anyone can edit.

[Introduction](#) • [Project Chat](#) • [Community Portal](#) • [Help](#)

Want to help translate? [Translate the missing messages.](#)

Welcome!

Wikidata is a free and open knowledge base that can be read and edited by both humans and machines.

Wikidata acts as central storage for the **structured data** of its Wikimedia sister projects including Wikipedia, Wikivoyage, Wiktionary, Wikisource, and others.

Wikidata also provides support to many other sites and services beyond just Wikimedia projects! The content of Wikidata is available under a [free license](#), exported using standard formats, and can be interlinked to other open data sets on the linked data web.

Get involved

Learn about Wikidata

- What is Wikidata? Read the [Wikidata introduction](#).
- Explore Wikidata by looking at a featured showcase item for author [Douglas Adams \(Q42\)](#).
- Get started with Wikidata's [SPARQL query service](#).

Contribute to Wikidata

- Learn to edit Wikidata: follow the [tutorials](#).
- Work with other volunteers on a subject that interests you: [join a WikiProject](#).
- Individuals and organisations can also [donate data](#).

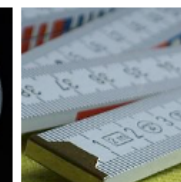
[Meet the Wikidata community](#)

Learn about data

New to the wonderful world of data? [Develop and improve your data literacy](#) through content designed to get you up feeling comfortable with the fundamentals in no time.



Item: [Earth \(Q2\)](#)



Property: [highest point \(P610\)](#)



custom value: [Mount Everest \(Q513\)](#)

Current highlights

- [Avengers: Endgame \(Q23781155\)](#)
- [Red Kelly \(Q706565\)](#)
- [The Visitation \(Q63444435\)](#)
- [Terje Moe Gustavsen \(Q4569851\)](#)

24

+

На месте жесткой посадки
Шереметьево найдены
snob.ru

Main page

Community portal

Project chat

Create a new Item

Create a new Lexeme

Recent changes

Random item

Query Service

Nearby

Help

Donate

Print/export

Create a book

Download as PDF

Printable version

In other projects

[Wikimedia Commons](#)

[MediaWiki](#)

[Meta-Wiki](#)

[Wikispecies](#)

[Wikibooks](#)

[Wikimania](#)

[Wikinews](#)

[Wikipedia](#)

[Wikiquote](#)

[Wikisource](#)

[Wikiversity](#)

[Wikivoyage](#)

[Wiktionary](#)

Tools

[What links here](#)

[Related changes](#)

[Special pages](#)

[Permanent link](#)

[Page information](#)

[Wikidata item](#)

In Wikipedia



Помощник по запросам



+ Фильтр

это частный случай понятия

человек



+ Показать

гражданство



Ограничение



```
1 SELECT ?person ?label_zh ?label_ru ?country
2 WHERE
3 {
4   ?person wdt:P31 wd:Q5.
5
6   ?person rdfs:label ?label_zh filter (lang(?label_zh) = "zh").
7   ?person rdfs:label ?label_ru filter (lang(?label_ru) = "ru").
8   ?person wdt:P27 ?country .
9   FILTER (?country IN (wd:Q159)) .
10 }
11
```

14:13:15 GMT+3,
17 мар. 2019 г.Данные
обновлены

2 минуты назад



2327 результатов за 5114 мс

</> Код

Скачать

Ссылка

Search



person	label_zh	label_ru	country
wd:Q4864	尤金·卡巴斯基	Касперский, Евгений Валентинович	wd:Q159
wd:Q162646	瓦迪姆·博羅夫謝夫	Вадим Владимирович Бровцев	wd:Q159
wd:Q163256	瓦莲京娜·古尼娜	Гунина, Валентина Евгеньевна	wd:Q159
wd:Q164309	法祖·阿里耶娃	Фазу Гамзатовна Алиева	wd:Q159
wd:Q122425	安德烈·萨金塞夫	Андрей Петрович Звягинцев	wd:Q159
wd:Q166498	維薩里翁	Тороп, Сергей Анатольевич	wd:Q159
wd:Q167113	卡倫·沙赫納扎羅夫	Карен Шахназаров	wd:Q159

DATASET

Parameters:

Categories:

- Country
- Type of an object

Words:

- Russian proper names
- Wikidata transliteration
- Xinhua-based transliteration

Metrics:

- Absolute Levenstein distance
- Normalized Levenstein distance
- Jaccard index

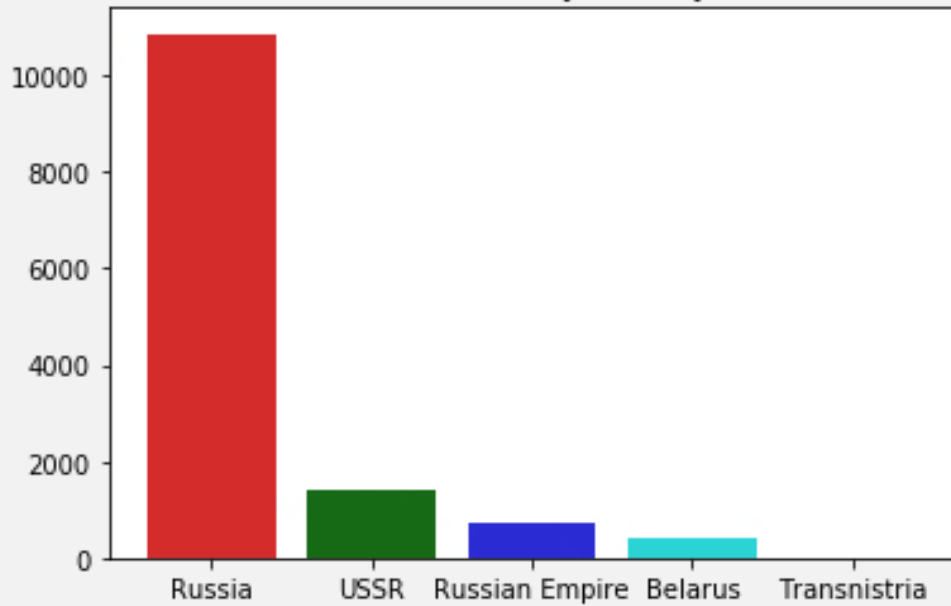
Country sample:

- Russian Federation
- Transnistria (Приднестровье)
- The USSR
- The Russian Empire
- Belarus
- Ukraine – pilot study

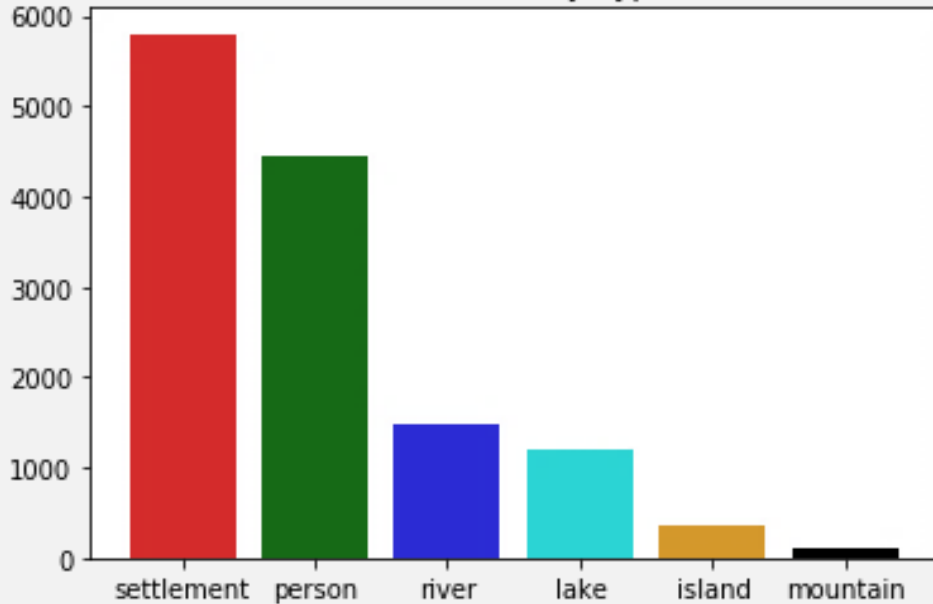
Type sample:

- Lakes
- Rivers
- Islands
- Mountains
- Settlements
- Person names

Distribution by Country



Distribution by Type



DATASET: OVERVIEW

13410 objects

Countries: 81% - Russia, 10% - USSR

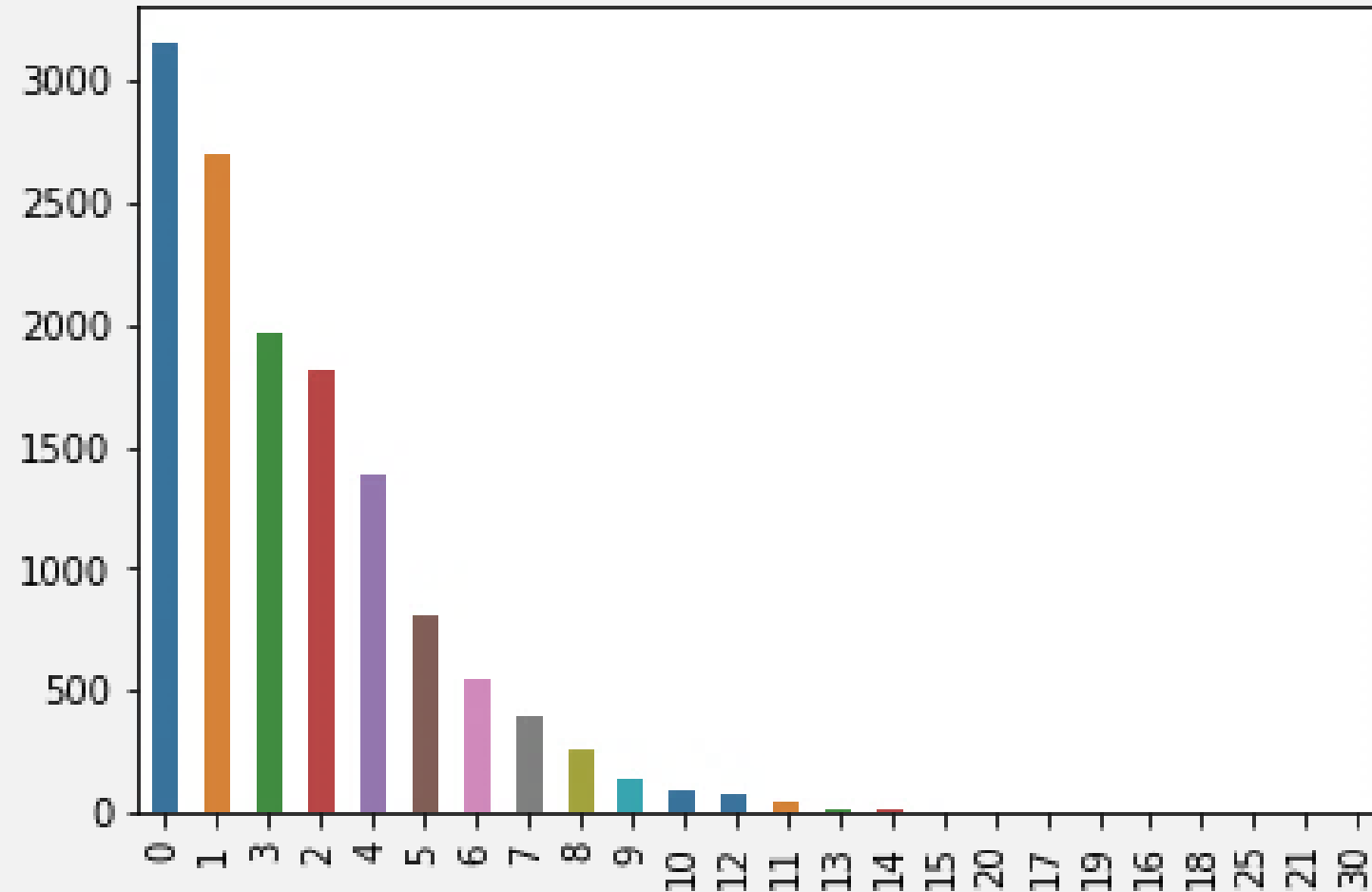
Types: 43% - settlements, 33% - names of persons, 11% - rivers

DATASET: STUDY

Levenshtein distance:

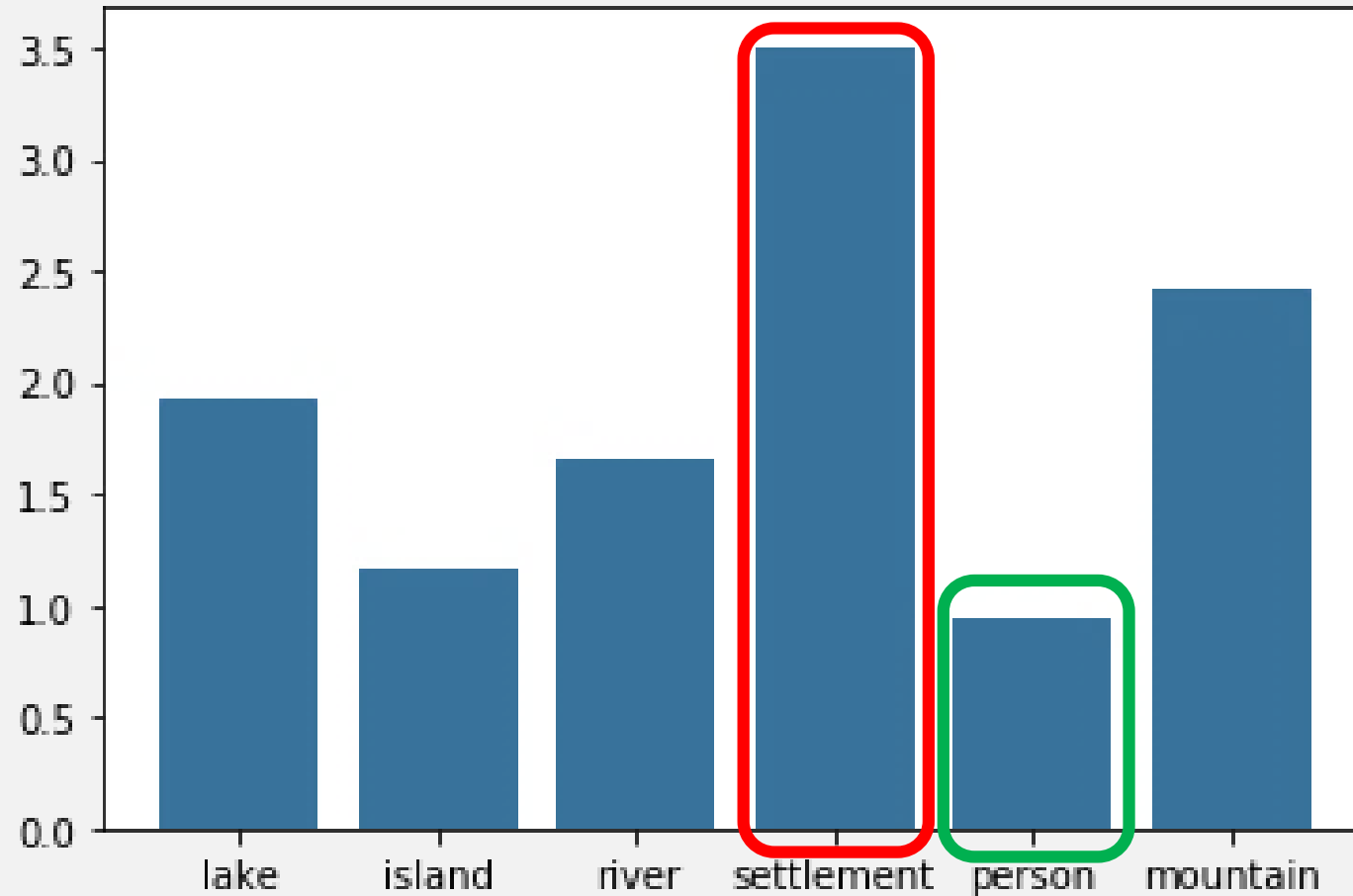
Absolute: $[0, \infty]$; 0 – identical strings,

$l - \infty$ - number of symbols to change



DATASET: STUDY

Mean Absolute Levenshtein Distance by Type



Classifiers?..

DATASET: RESULTS

Xinhua rules are not so bad for personal names...

...but are hardly applicable to other types of objects

There might be some semantic elements which are translated to Chinese...

3. SEMANTIC ELEMENTS IN TRANSLITERATION

CLASSIFIERS: OVERVIEW

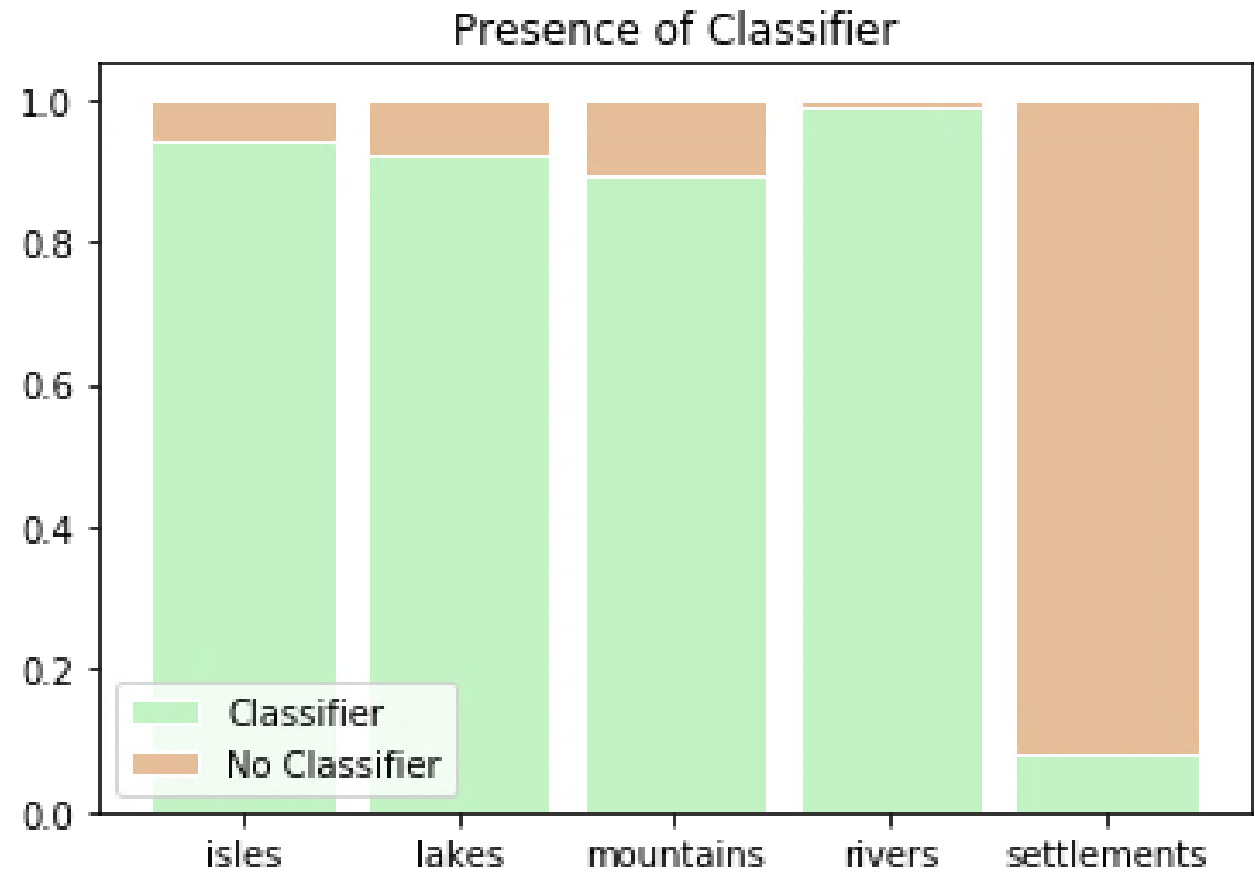
One syllable

In the end of the word

No classifiers for personal
names

- 伏尔加河 fú'ěrjiā hé
- 奥涅加湖 àonièjiā hú
- 萨哈林岛 sàhālín dǎo
- 弗拉基米尔市 fúlājīmǐ'ěr shì

CLASSIFIERS: FREQUENCY



IRRESISTIBLE
COMPULSION TO
SEMANTIC ELEMENTS

Novi Zeland VS New York

Новая Зеландия VS
Нью-Йорк

How strong is it in
Chinese?

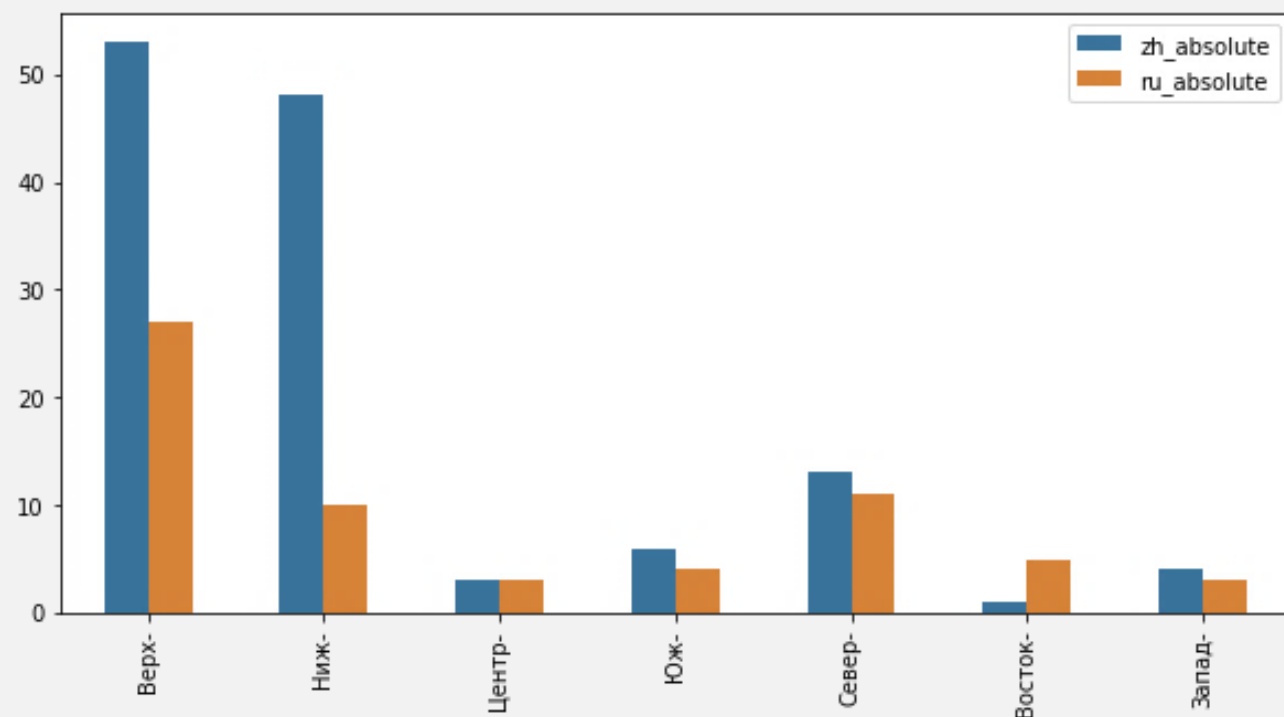
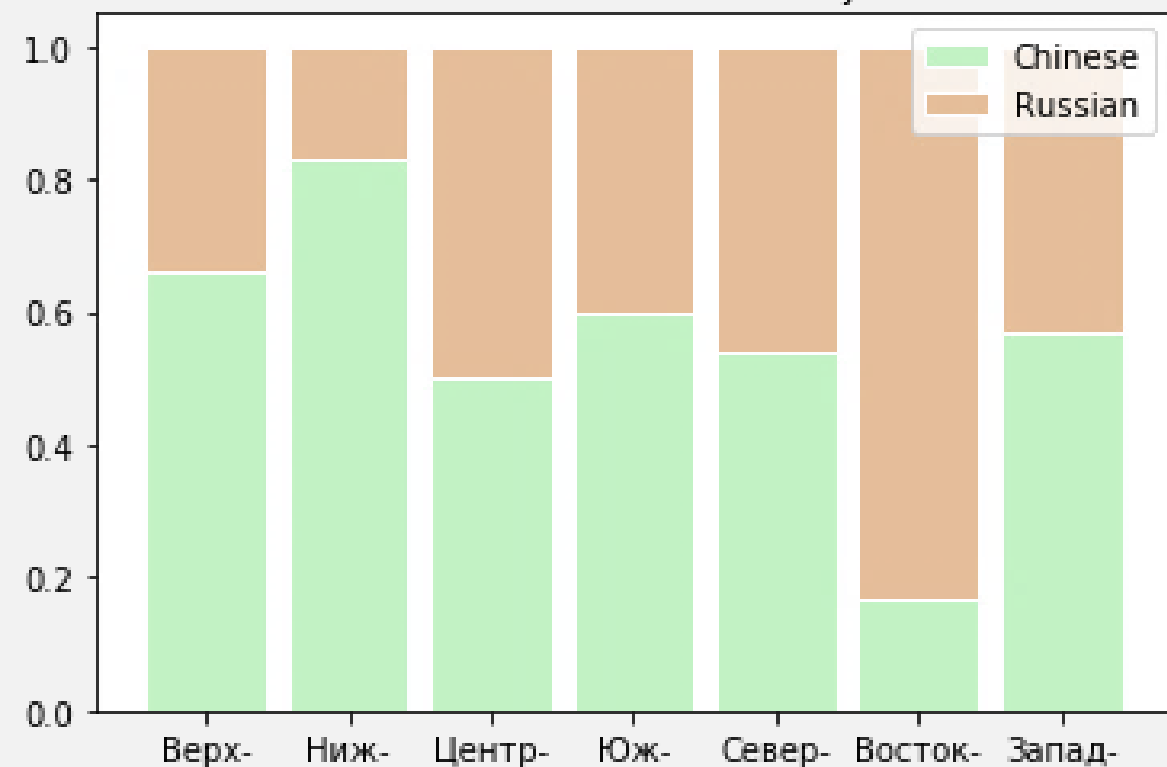
PARTIAL SEMANTIC TRANSLATION

Based on the rough data overview, we split the possible variants into 3 groups:

- “Spatial” Adjectives (“Higher”, “Southern”, “Central” etc.)
- Other Adjectives (“Greater”, “New”, “Red” etc.)
- Affixes (“Trans-”, “-upon-” etc.)

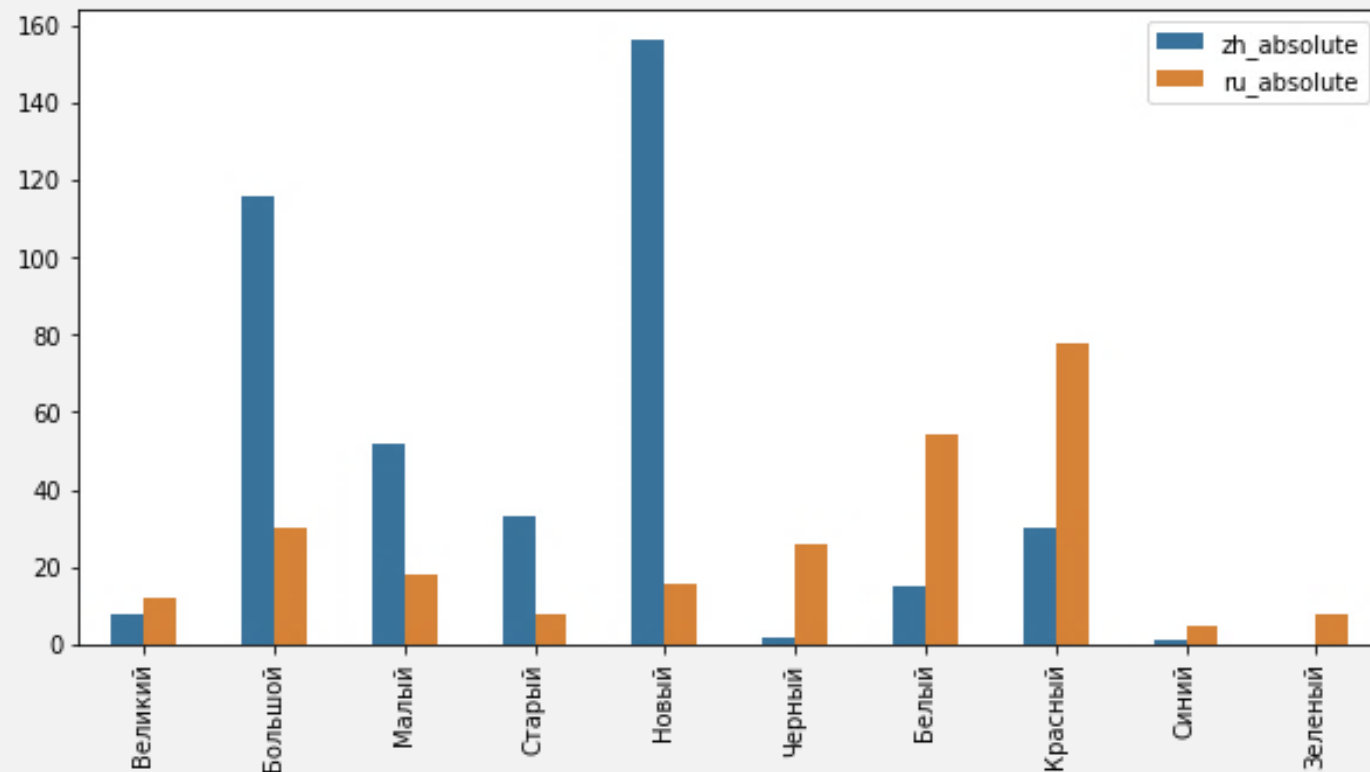
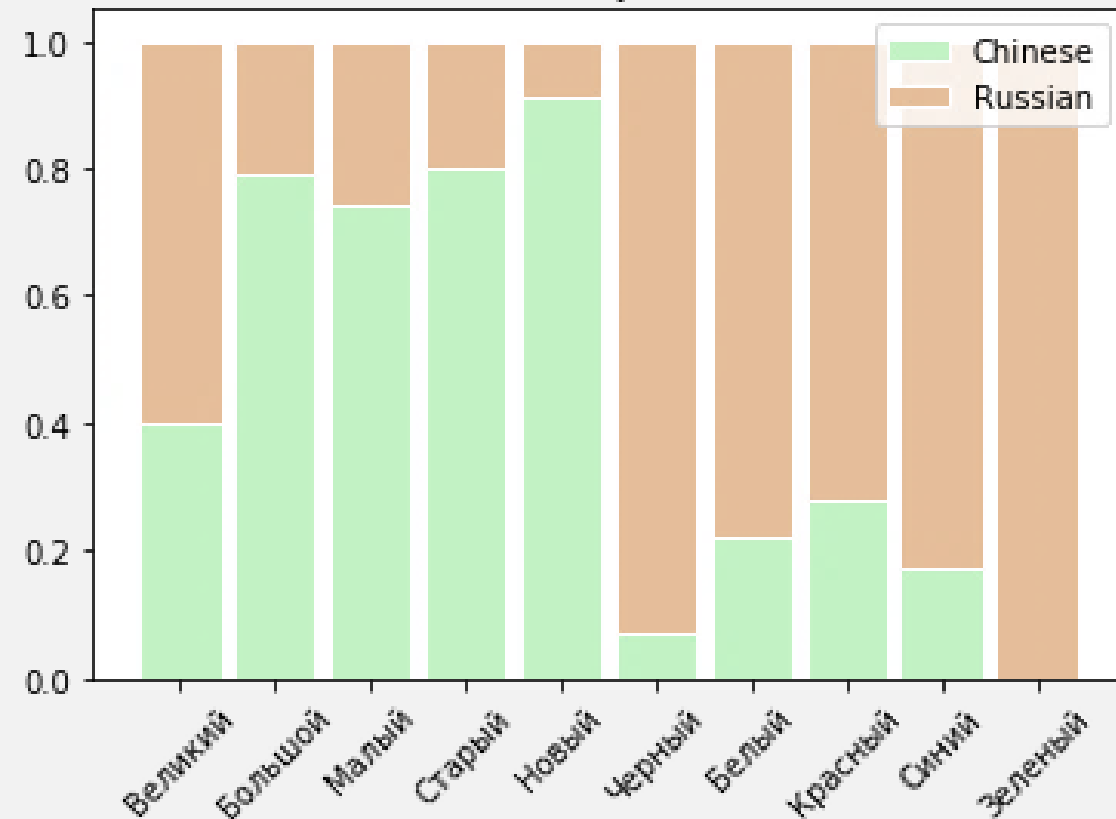
PARTIAL SEMANTIC TRANSLATION: “SPATIAL” ADJECTIVES

Orientational Prefixes and Adjectives



PARTIAL SEMANTIC TRANSLATION: OTHER ADJECTIVES

Other Adjectives



PARTIAL SEMANTIC TRANSLATION: AFFIXES

The least likely to be semantically transformed

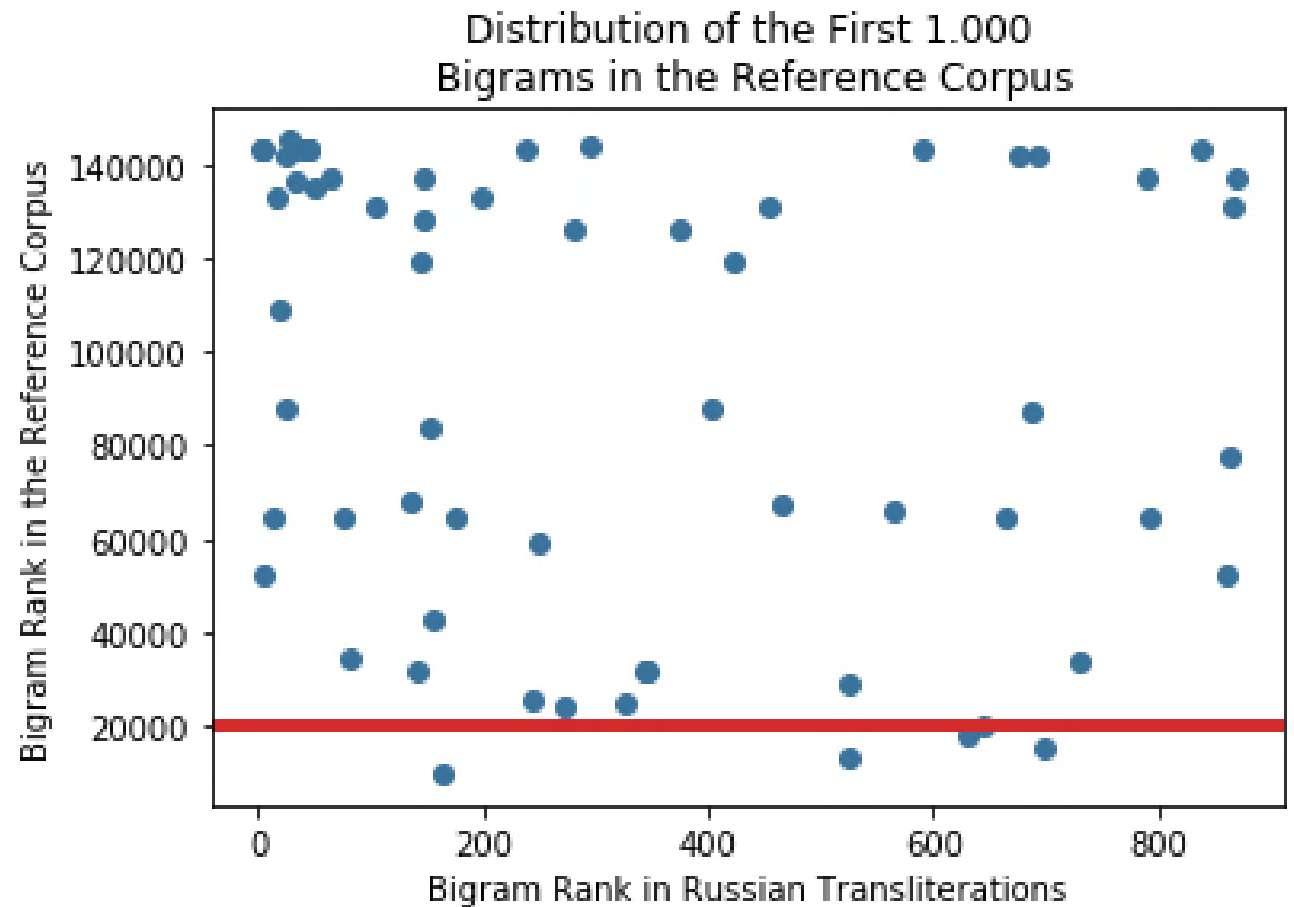
Two occurrences:

- Забайкальск > 后-贝加尔斯克 > hòu (behind) -bèijiā'ěrsīkè
- Ростов-на-Дону > 顿河畔罗斯托夫 > dùnhé pàn luósītōufū

4. N-GRAM FREQUENCY

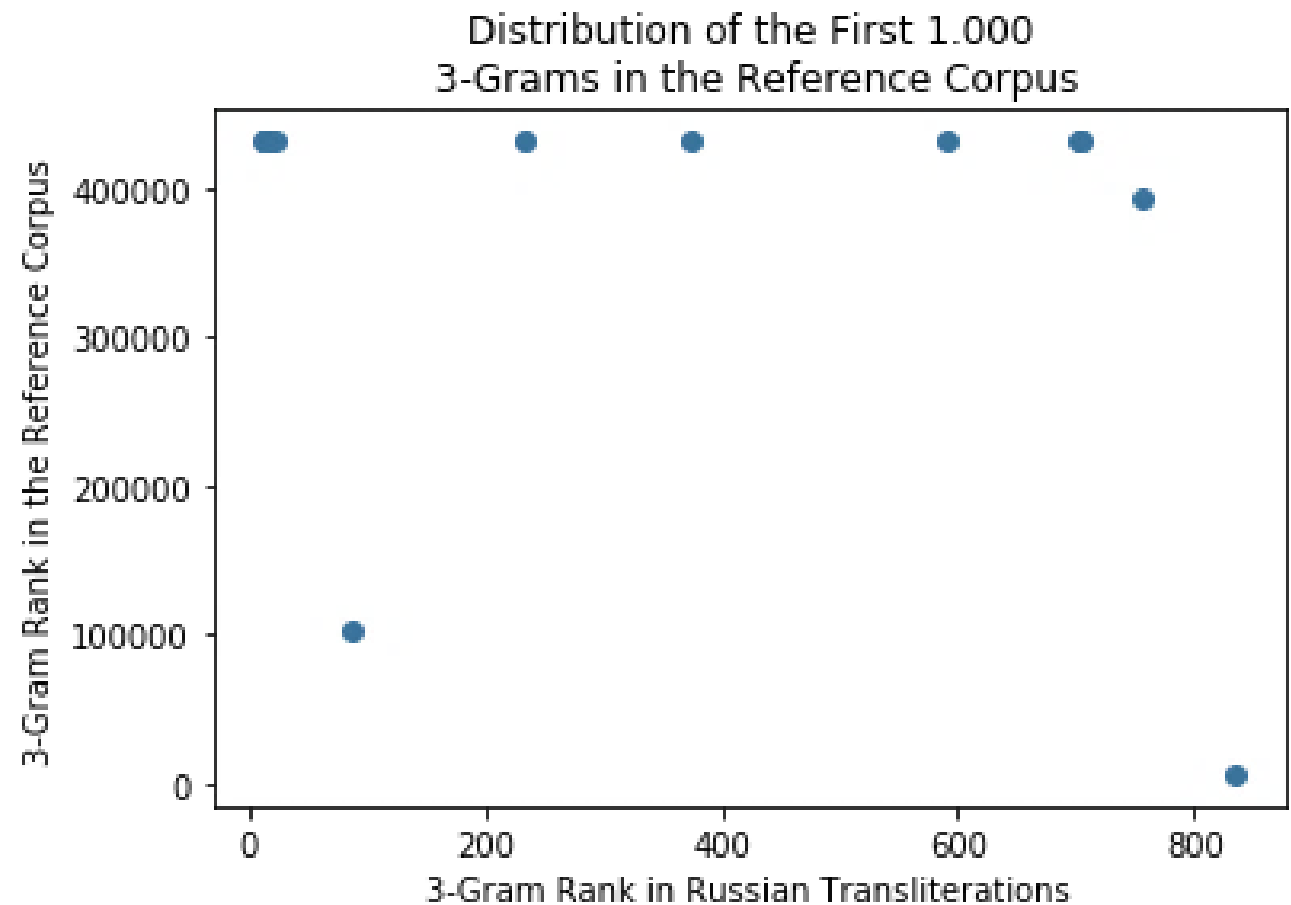
BIGRAMS

936 of Bigrams are not met in the reference corpus at all



3-GRAMS

988 of 3-Grams are not met in the reference corpus at all



RESULTS

- Generally the trends of Russian loanwords' phonetic adaptation are similar to those of other European languages
- But the official (Xinhua) prescriptions do not cover them
- The Chinese meaningful characters can be an efficient hint in limiting the loanwords
 - Adjective-like affixes – from the left
 - Classifiers (except personal names) – from the right
- The majority of N-grams in transliterations are rare in the usual Chinese texts

FUTURE PERSPECTIVES

Implementation of the statistics in Chinese NLP algorithms

- Hidden Markov models
- Seq2Seq Neural networks

Analysis of a bigger dataset

- BaiduPedia (百度百科) – bigger than Wikipedia in 5 European languages altogether
- More oriented on PRC

conduct the same research on other European languages

- Experts in these language are needed – we invite you to take part!

Thank you for your
attention!

Hvala na pažnji!

非常感谢！

Kirill Semenov,
HSE – Moscow
kir.semenow@yandex.ru

ADDITIONAL MATERIALS

PROBLEMS: ALL OVER RUSSIA

* 红场 - hóng chǎng – Red
Square

红肠 - hóng cháng – Red
Guts



PHONETIC ADAPTATION OF
THE RUSSIAN WORDS IN
CHINESE: THE OT APPROACH

- Sources:
 - 汉语外来词词典 (Chinese Loanword Dictionary), 1984 – 387 words
 - БКРС-Online (Big Chinese-Russian Dictionary Online) – 1494 items

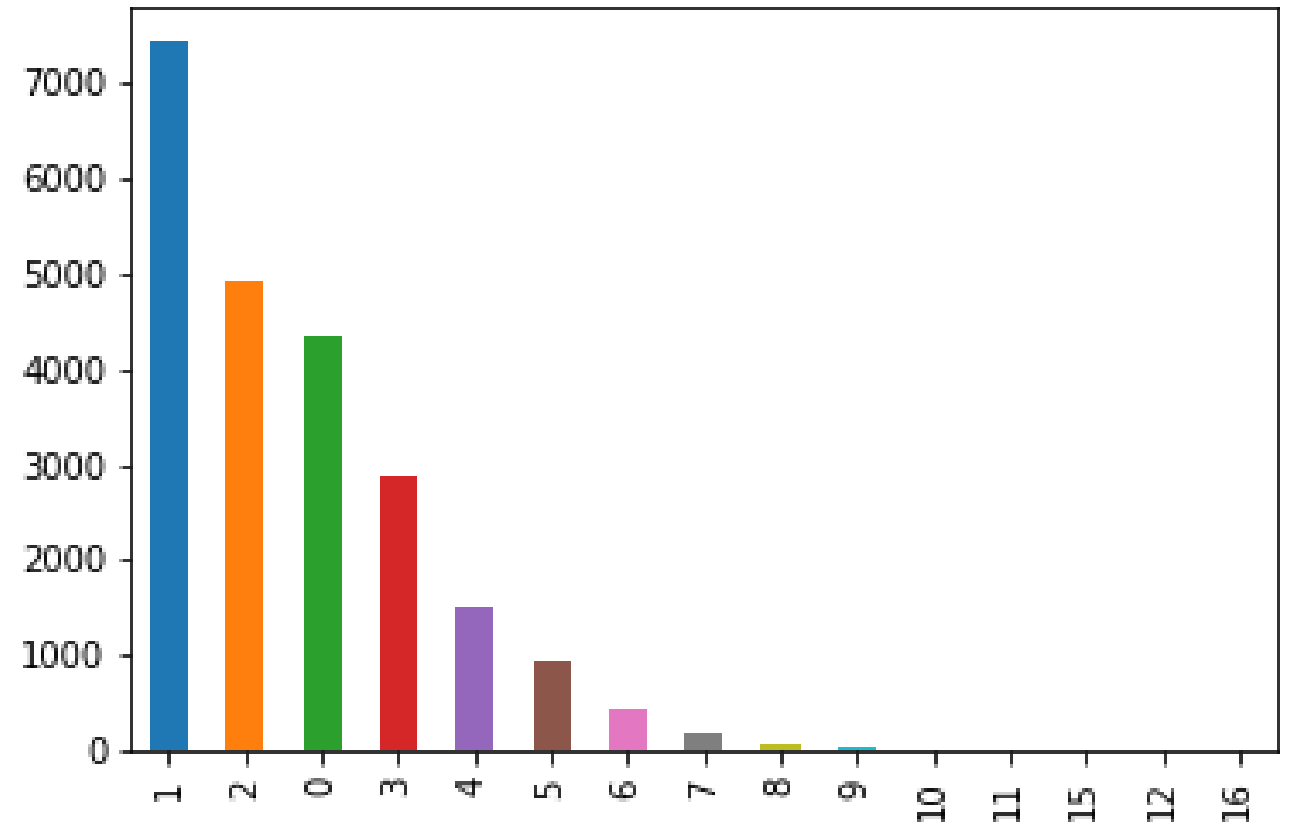
CASE STUDY: MATERIAL

Dictionaries and grammar of
Russian-Chinese pidgin
(Perekhval'skaya 2008)

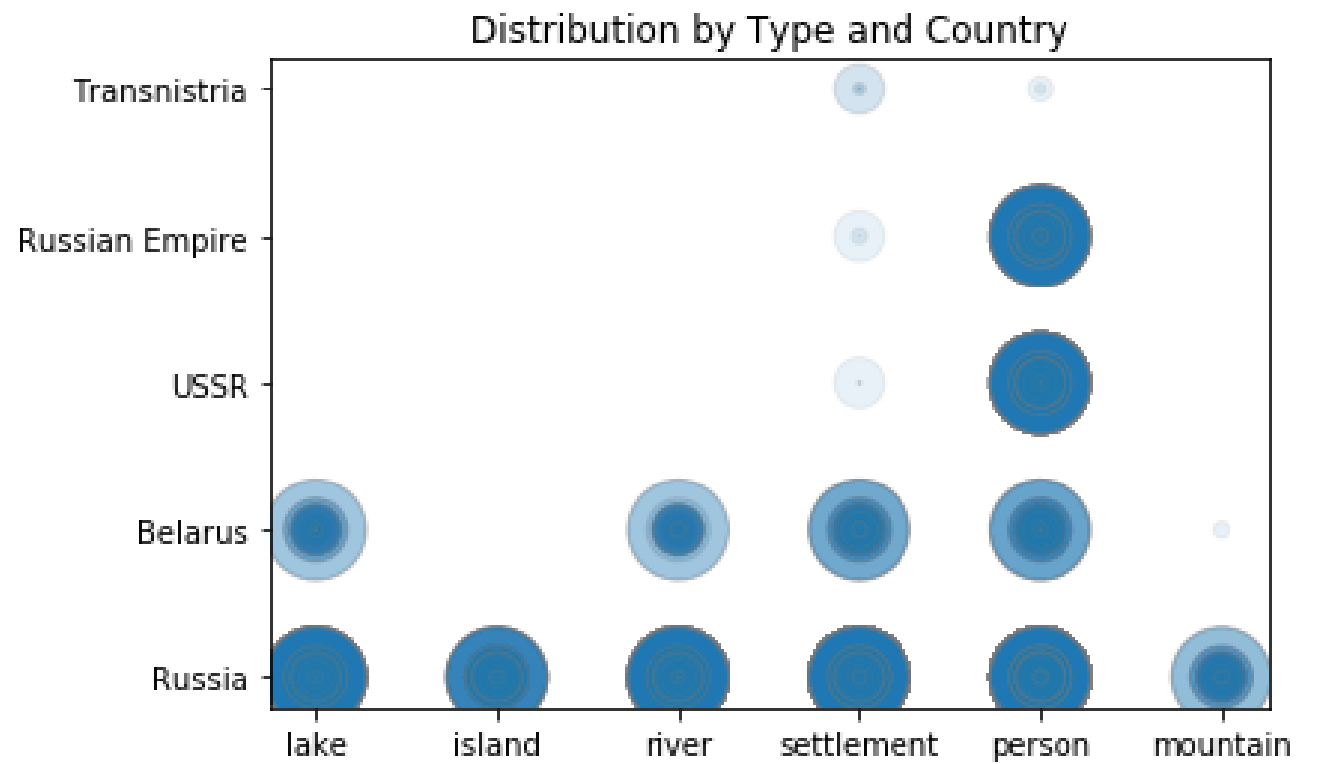
Grammar of the Northern
Chinese dialects (Zavyalova 1996)

Dictionaries of the Northern
Chinese dialects and of the
modern Chinese

LEVENSTEIN
DISTANCE FOR
UKRAINIAN
SETTLEMENTS



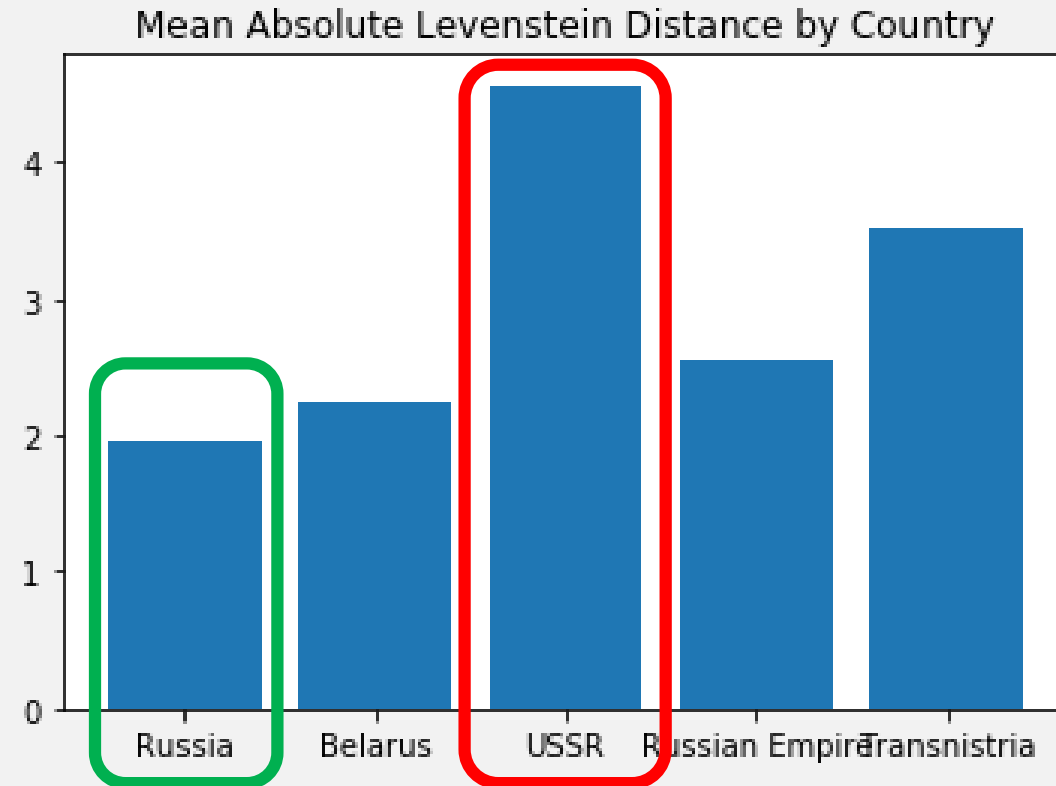
DATASET: OVERVIEW



DATASET: STUDY

Character diversity in Wikidata items	877
Character diversity of Xinhua-based transliterations	260

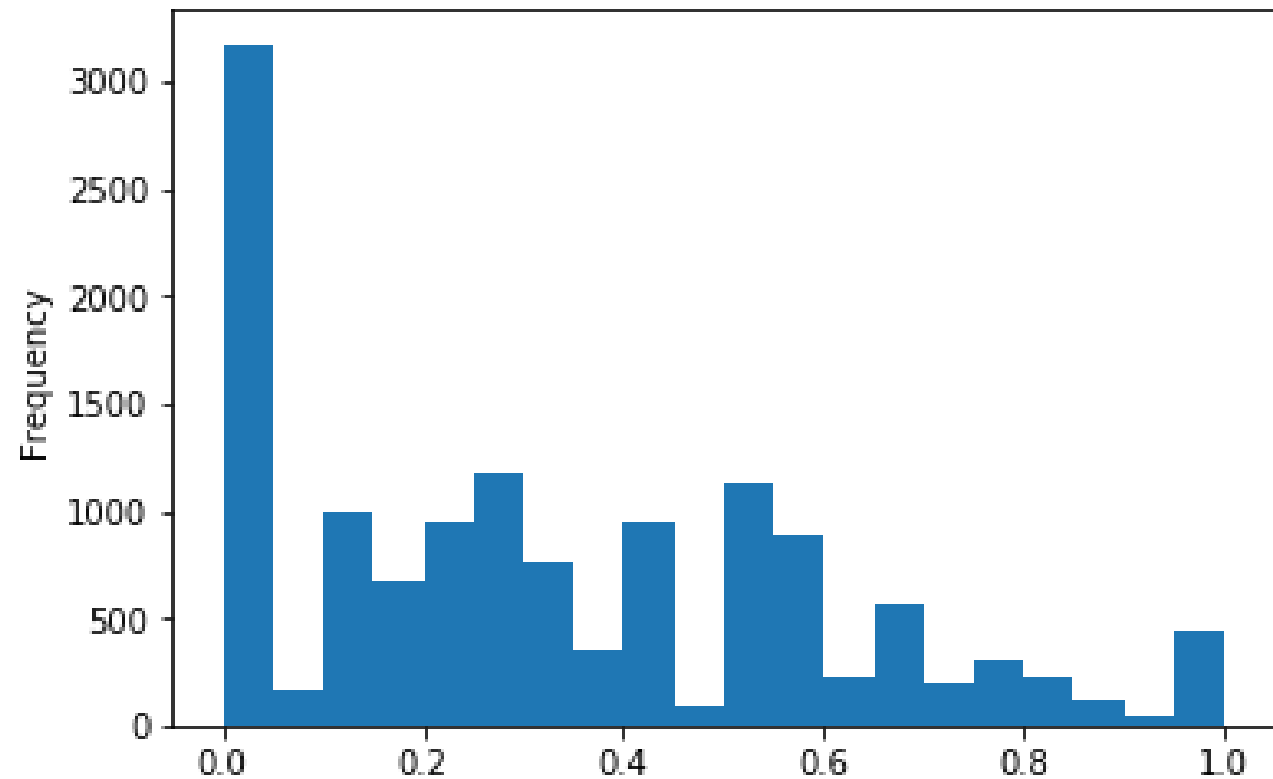
DATASET: DISTRIBUTION BY TYPE



Communist names?..

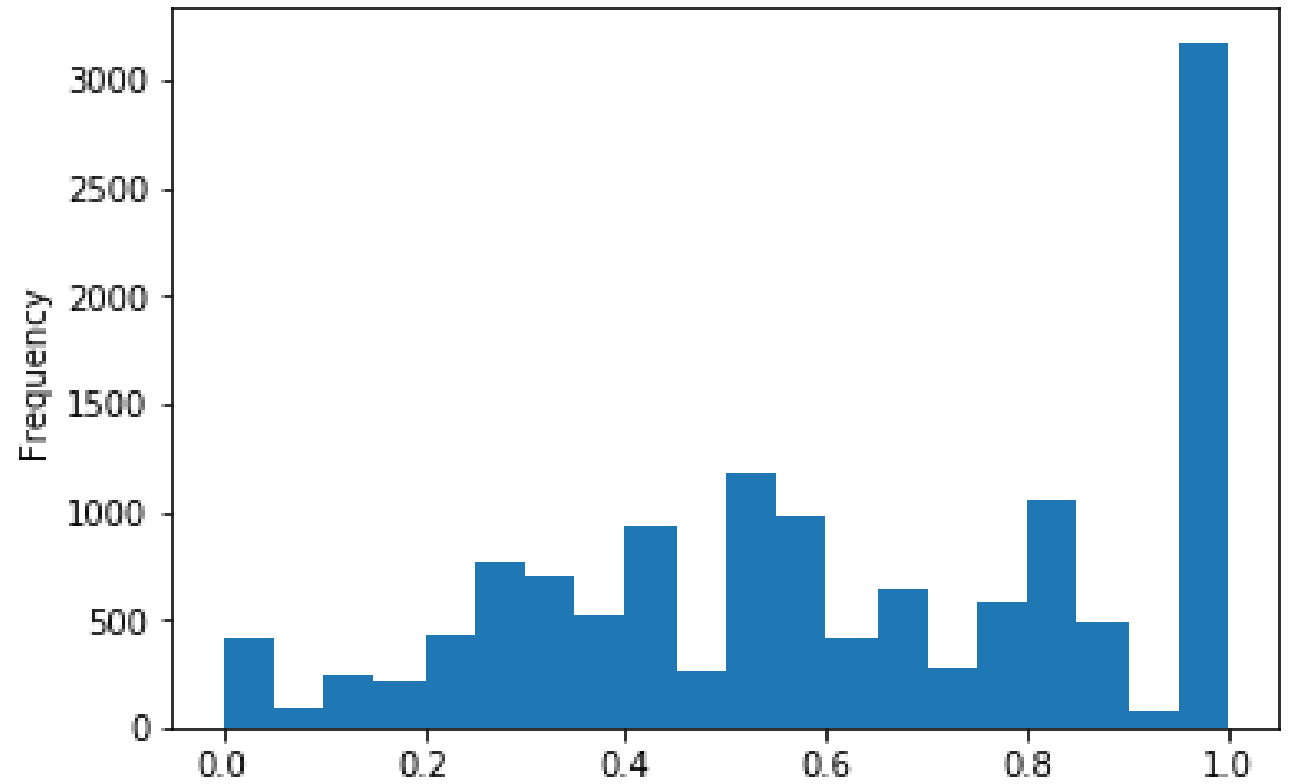
DATASET: STUDY

- Normalized Levenstein distance:
- $[0, 1]$; 0 – identical strings, 1 – totally different strings

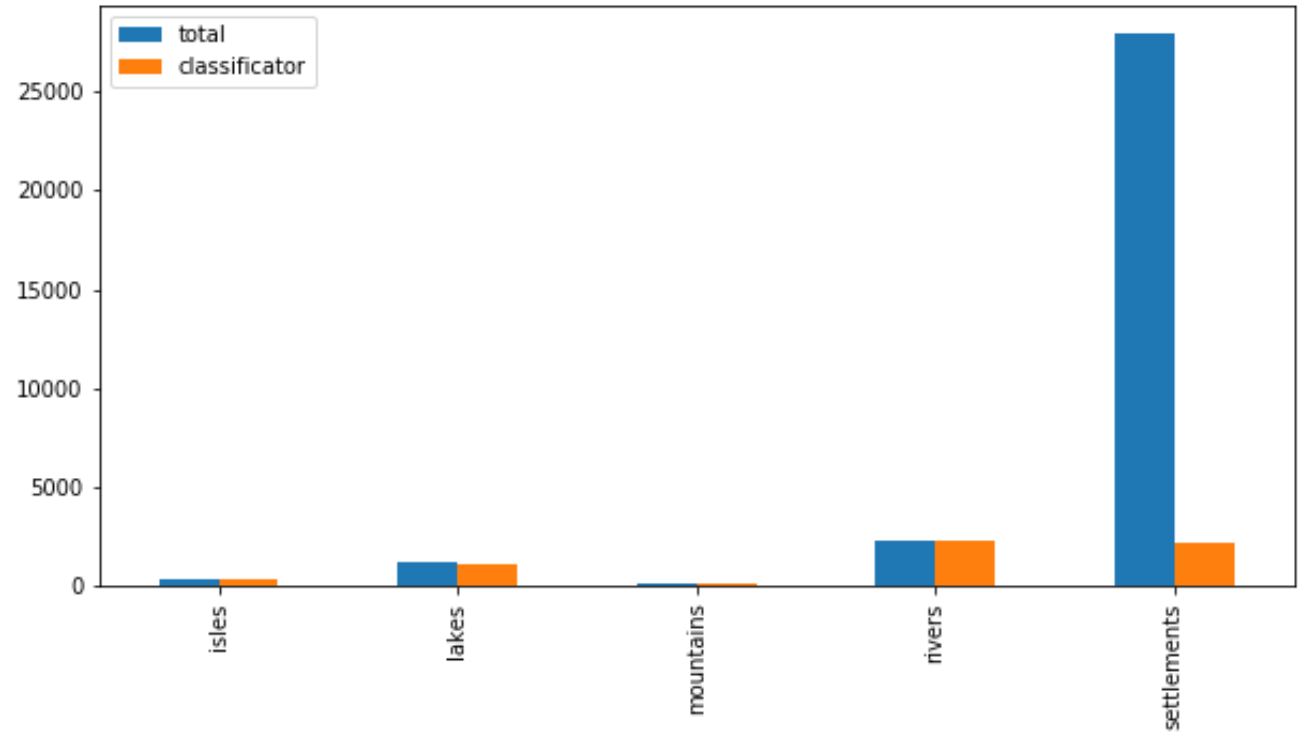


DATASET: STUDY

Jaccard Index: $[0, 1]$; 0 –
the same set of
symbols, 1 – no overlap



CLASSIFIERS: ABSOLUTE NUMBERS



TOTAL SEMANTIC TRANSLATIONS

- East Asian Toponyms
 - Especially former Japanese territories – Kuril islands, Sakhalin, etc.
 - On-border settlements:
 - 海参崴 hǎishēnwǎi (Sea cucumber river bend) = Владивосток
- Communist Toponyms
 - Октябрьское, Первомайское, etc.
- Other literal translations:
 - Белая (деревня)
 - Старое (село)
 - Аэропорт (район)
 - Чистые пруды

N-GRAMS: CHINESE WIKIDATA 3-GRAMS

维奇·739

夫卡\$ 491

斯基\$ 466

夫斯\$ 434

科耶\$ 410

斯克\$ 396

亚历山 355

斯科耶 350

\$亚历 280

科夫\$ 273

历山大 253

诺夫\$ 251

耶维奇 242

山大·231

夫斯基 215

尼古拉 210

米哈伊 200

耶夫\$ 189

谢尔盖 185

\$谢尔 182

诺耶\$ 164

\$弗拉 160

洛夫\$ 157

诺维奇 154

\$尼古 148

拉基米 146

尔盖·146

弗拉基 145

罗维奇 142

米尔·137

N-GRAMS: RUSSIAN WIKIDATA 4-GRAMS

кий\$ 2550

ский 2508

кое\$ 2204

ское 2100

овск 1385

вско 937

инск 904

вски 864

нски 834

вич 753

нско 655

ович 583

андр 536

лекс 525

алек 521

\$але 421

евск 408

ксан 390

санд 379

нико 379

екса 376

енск 344

ий\$к 339

евич 330

анов 319

ова\$ 318

ков\$ 300

новс 295

ковс 285

ое\$к 279

ALGORITHM APPLICATION: PART I

- Chinese Loanwords' Dictionary (汉语外来词词典), Shanghai, 1984
- 378 words of Russian origin
 - Or borrowed into Chinese via Russian

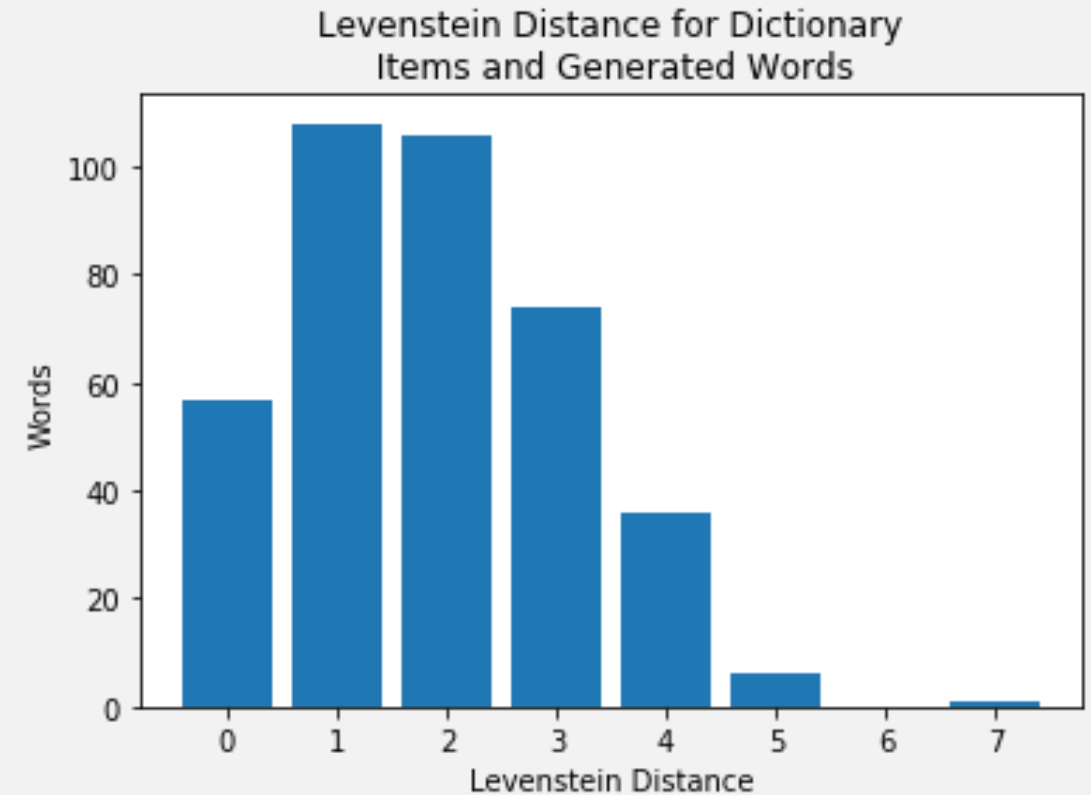
ALGORITHM APPLICATION: PART I

- Dataset:
 - Chinese dictionary occurrence
 - Russian analogue
 - Xinhua-generated transliteration of the Russian word
 - Levenstein distance between the dictionary and Xinhua words

id	word	origin	Xinhua	levenstein_word
0	阿尔申	аршин	阿尔申	0
1	阿尔西非	альсифер	阿利西费尔	3
2	阿法林	афалина	阿法莉娜	2
3	阿留米特	алюмит	阿柳米特	1
4	阿依尔	аил	艾尔	2
5	阿札林	азарин	阿扎林	1
6	艾费勃	офеб	奥费布	2
7	艾米利通	эмиритон	埃米里托恩	4
8	艾木兴	эмшер	埃姆舍尔	4
9	艾特纳	этноэ	埃特内	2
10	艾匹配	эпипэ	埃皮佩	3
11	爱特罗尔	этрол	埃特罗尔	1
12	安诺	анау	阿瑙	2

ALGORITHM APPLICATION: PART II

Character diversity in dictionary items	319
Character diversity of Xinhua-based transliterations	168



- Explanation: there is a bigger variance of the Chinese characters in the dictionary, and the problem is in multiple choice of a character based on one phonetic reading:
 - “li”: 利 (Xinhua)
利, 里, 理, 立, 列 (dictionary)