

# Частотный словарь параллельного корпуса

## **Студенты:**

София Землянская,  
Аполлинария Карпова,  
Александра Кузьмина,  
Таисия Лукьянова,  
Александра Нужненко,  
Анастасия Сычева,  
Анастасия Фирсова,  
Елизавета Шемшурина

## **Руководитель:**

Ольга Николаевна Ляшевская

# Цель

- Создать базу данных, основу для частотного словаря, в котором представлены **переводы** слова с одного языка на другой и **частотность** этих переводов. Информация о частотности основана на данных **параллельных корпусов**.

# Для чего пригодится словарь

- **Лексическая типология:** теоретический анализ фреймов, подбор контекстов для опросов;
- **Психолингвистика:** при изучении билингвизма создаются эксперименты, в которых представлены одни и те же по смыслу предложения или слова на двух языках; словарь поможет подобрать наиболее подходящий перевод;
- **Преподавание иностранного языка;**
- **Машинный перевод:** на основе словаря можно создать датасет для оценки машинного перевода

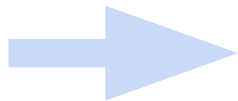
# Данные

- **Русско-английский и англо-русский корпуса**, в которых реализовано выравнивание по предложениям
- **Задача:** подготовить пословное выравнивание и частотные списки пар слов → пословное выравнивание обоих корпусов с помощью скрипта `fast_align`

She went to the shop.



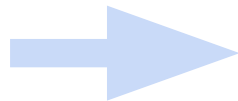
Она пошла в магазин.



She went to the shop.



Она пошла в магазин.



eng	rus	num
she	она	n
go	пойти	m
to	в	k
shop	магазин	j

# Обработка данных

- Предварительная обработка данных: преобразование файлов xml в txt формат, который поддерживает text\_align: eng\_text ||| rus\_text или rus\_text ||| eng\_text;
- Апробация двух способов: сначала лемматизация, потом выравнивание и сначала выравнивание, потом лемматизация;
- Выбран вариант лемматизация - выравнивание, так как в нем было меньше мусора и больше адекватных соответствий (например, нет случаев, когда английский артикль приравнивается к русскому слову, так как артикли удаляются заранее)
- Добавление частей речи после слов;
- Английский – с помощью NLTK, русский – с помощью rymorphy2, теги nltk переделаны в rymorphy2 для однородности данных.

# fast\_align (C++, Cmake)

## Формат ВХОДНЫХ ДАННЫХ:

doch jetzt ist der Held gefallen . ||| but now the hero has fallen .  
neue Modelle werden erprobt . ||| new models are being tested .  
doch fehlen uns neue Ressourcen . ||| but we lack new resources .

## Alignment:

- source–target (left language–right language) alignments
- target–source (right language–left language) alignments
- atools: симметрия с использованием различных стандартных эвристик

# fast\_align (C++, Cmake)

## Формат ВЫХОДНЫХ ДАННЫХ:

- “Pharaon format” (i-j)

0-0 1-1 2-4 3-2 4-3 5-5 6-6

0-0 1-1 2-2 2-3 3-4 4-5

0-0 1-2 2-1 3-3 4-4 5-5

## Создание списка слов:

8     0.0909    генрих    henry

1     0.0114    генрих    martin

# Сложности выравнивания

- Неправильное соотнесение слов (например, в пару к самостоятельным словам подобраны служебные)
- Самый частотный перевод  $\neq$  самый правильный

id	count	probability	english	russian
0	38	0.0219	bird	птица
1	19	0.0109	bird	птичий
2	71	0.0408	bird	и



# Результат

Русско-английский корпус наречий - 11730 примеров;

Англо-русский корпус наречий- 22230 примеров;

Разметка результата: 0 - перевод неверный , 1 - перевод верный

index	count	probability	russian	english	result
25572	10	0.0075	трудно_ADVB	life_NOUN	0
25573	3	0.0022	трудно_ADVB	hard_ADVB	1
25574	3	0.0022	трудно_ADVB	began_VERB	0
25575	3	0.0022	трудно_ADVB	time_NOUN	0
25576	3	0.0022	трудно_ADVB	since_PREP	0
...			...		

Русско-английский корпус: 1328 верных и 10402 неверных переводов

# Сложности разметки

- Разные грамматические пометы у одних и тех же слов и слов, совпадающих по значению; соответствия разных частей речи.

index	count	probability	russian	english
86236	29	0.0122	тихо_ADVB	<b>quiet_ADJF</b>
86263	4	0.0017	тихо_ADVB	<b>quiet_NOUN</b>

- Случаи, в которых перевод — элемент соответствия;

index	count	probability	russian	english
103802	19	0.1	исподлобья_ADVB	looked_VERB
103806	3	0.0158	исподлобья_ADVB	beneath_PREP

**исподлобья** нареч.  
общ. frowningly; from under the brows (distrustfully, sullenly)  
from under one's eyebrows (He looked up at me from under his eyebrows  
a sullen look (в контексте: he gave me a sullen look

- Обратные случаи: оригинал — элемент соответствия.

index	count	probability	russian	english
106201	3	0.0057	пристально_ADVB	stare_VERB

# Контексты

Как составлялась таблица:

- Для каждой пары из базы данных - поиск пар предложений, где есть оба слова

Но:

- Не всегда точные соответствия частотностям из базы
- Пока сделана только для русско-английского корпуса

*Частотный словарь параллельного корпуса — это коллекция текстов и их переводов на разных языках. Словарь поможет исследовать различия в активности лексики двух языков.*

Выберите язык-источник



Выберите язык перевода



Введите слово



Введите слово



Искать

*Частотный словарь параллельного корпуса — это коллекция текстов и их переводов на разных языках. Словарь поможет исследовать различия в активности лексики двух языков.*

Русский

?

Английский

?

Введите слово

?

Введите слово

?

Искать

*Частотный словарь параллельного корпуса — это коллекция текстов и их переводов на разных языках. Словарь поможет исследовать различия в активности лексики двух языков.*

?

?

?

?

машина — car — 1056 hits

Скачать результат запроса

**Машина** принадлежит Мелиссе Джой Блэк из Мичигана.

Uh, **car** belongs to a Melissa Joy Black from Michigan.

?

Мы были женаты несколько лет, и у нас было двое детей и три **машины**, но их он оставил себе.

We were married for several years and had two children and three **cars**, and he kept the cars.

?

⋮

⋮

- И мы так думали, но **машина** не застрахована.

"So we thought, but the **car's** not insured."

?



машина — car — 1056 hits

Скачать результат запроса

**Машина** принадлежала  
Мичигана.

Мы были женаты,  
было двое детей,  
оставил себе.

- И мы так дума  
застрахована.

Сначала он размышлял о чем-то своем, затем медленно произнес:  
- И мы так думали, но **машина** не застрахована.  
Повисло напряженное молчание.

**Автор:** Василий Петров  
**Название:** "Происшествие"

---

At first he was thinking about something of his own, then he said slowly:  
"So we thought, but the **car's** not insured."  
There was a tense silence.

**Переводчик:** Christopher R. Wood  
**Название:** The incident

back from ?

and had two  
pt the cars. ?

insured." ?

*Частотный словарь параллельного корпуса — это коллекция текстов и их переводов на разных языках. Словарь поможет исследовать различия в активности лексики двух языков.*

?

?

?

?

Русский

**Машина** принадлежит Мелиссе  
Джой Блэк из Мичигана.

К полуночи **машина** хорошо  
работала и полным ходом.

⋮

— Бронированная **машина** номер  
три. Мне необходимо срочно  
найти ее.

Английский

Uh, **car** belongs to a Melissa Joy Black  
from Michigan.

By midnight the **machinery** of his  
command was working thoroughly  
and efficiently.

⋮

"Armored **truck** number three. I need  
to find it."

Частота

?

1056 hits

>>>

?

550 hits

>>>

?

35 hits

>>>

машина — car — 1056 hits

Скачать результат запроса

**Машина** принадлежит Мелиссе Джой Блэк из Мичигана.

Uh, **car** belongs to a Melissa Joy Black from Michigan.

?

Мы были женаты несколько лет, и у нас было двое детей и три **машины**, но их он оставил себе.

We were married for several years and had two children and three **cars**, and he kept the cars.

?

⋮

⋮

- И мы так думали, но **машина** не застрахована.

"So we thought, but the **car's** not insured."

?

*Частотный словарь параллельного корпуса — это коллекция текстов и их переводов на разных языках. Словарь поможет исследовать различия в активности лексики двух языков.*



Русский

Она взяла напрокат **машину** и направилась к Лонг-Айленду. Она направлялась совершить ограбление.

А от **автомобиля** до космических полетов - всего лишь несколько десятилетий.

⋮

Они отправились в путешествие в личном **вагоне** Джейми.

Английский

She rented a **car** and headed for Long Island. She was on her way to commit a burglary.

Yet only decades from the **car** into space.

⋮

They made the journey in Jamie's private railway **car**.

Частота

?

1056 hits

>>>

?

550 hits

>>>

?

35 hits

>>>

машина — car — 1056 hits

Скачать результат запроса

**Машина** принадлежит Мелиссе Джой Блэк из Мичигана.

Uh, **car** belongs to a Melissa Joy Black from Michigan.

?

Мы были женаты несколько лет, и у нас было двое детей и три **машины**, но их он оставил себе.

We were married for several years and had two children and three **cars**, and he kept the cars.

?

⋮

⋮

- И мы так думали, но **машина** не застрахована.

"So we thought, but the **car's** not insured."

?





# Перспективы проекта

- доделать сайт и выложить в общий доступ
- добавить другие пары языков:
  - русский-немецкий, русский-французский, русский-испанский и пр.
  - без русского: английский-немецкий, французский-испанский и пр.
- сохранять метайнформацию (автор, название произведения, год и др.)
- добавить тексты разных жанров
- снабдить некоторые из слов тэгами *разг.*, *вульг.*, *уст.* и пр.

# Наша команда

- Аполлинария Карпова: автоматическое выравнивание;
- Анастасия Фирсова: автоматическое выравнивание;
- София Землянская: чистка и анализ данных;
- Александра Кузьмина: чистка и анализ данных;
- Таисия Лукьянова: чистка и анализ данных;
- Анастасия Сычева: чистка и анализ данных;
- Елизавета Шемшурина: чистка и анализ данных, подготовка контекстов для списков слово-перевод;
- Александра Нужненко: дизайн сайта.

Спасибо за внимание!