

Параллельные корпуса

Темы курса

- Различные параллельные корпуса в составе НКРЯ (сложности создания)
- Alignment и как его делать
- Разметка
- Методы анализа, исследования на материале корпусов
- Построение датасета на примере пассива
- Извлечение контекстов из корпуса
- Семантические карты
- Графы
- Case study: Проект по церковно-славянским текстам
- Case study: Проект по текстам Нового Завета

Темы курса

- Различные параллельные корпуса в составе НКРЯ (сложности создания)
- Alignment и как его делать
- Разметка
- Методы анализа, исследования на материале корпусов
- Построение датасета на примере пассива
- Извлечение контекстов из корпуса
- Семантические карты
- Графы
- Case study: Проект по церковно-славянским текстам
- Case study: Проект по текстам Нового Завета
- Ваши исследования!

Базовые понятия

Устройство

Параллельный корпус представляет собой коллекцию оригинальных текстов на языке L_1 с их переводами на один или более языков $L_2...L_n$, важным атрибутом такого корпуса является наличие в нем **выравнивания**, то есть наличие установленных соответствий между текстовыми единицами.

Устройство

Параллельный корпус представляет собой коллекцию оригинальных текстов на языке L_1 с их переводами на один или более языков $L_2...L_n$, важным атрибутом такого корпуса является наличие в нем **выравнивания**, то есть наличие установленных соответствий между текстовыми единицами.

Как правило тексты в параллельном корпусе выровнены по предложениям (реже по фрагментам текста, абзацам), иногда в корпусе может быть также и пословное выравнивание, когда соответствия устанавливаются между словами.

Переводные единицы

Предложение (или несколько предложений) на языке оригинала с соответствующими ему переводами на представленные в корпусе языки называется **переводной единицей**.

EN: **He** **was interrupted** by a knock on the door.

RU: В дверь постучали.

DE: Ein Klopfen an der Tür **unterbrach** ihn.

ITA: **Fu interrotto** da qualcuno che bussava alla porta.

ES: **Le interrumpieron** unos golpes en la puerta.

CZ: **Přerušilo ho** zaklepaní na dveře.

BG: **Прекъсна го** почукване на вратата.

SE: **Han avbröts** av en knackning på dörren.

FR: **Il fut interrompu** par des coups frappés à la porte.

Переводные единицы

Предложение (или несколько предложений) на языке оригинала с соответствующими ему переводами на представленные в корпусе языки называется **переводной единицей**.

EN: **He was interrupted** by a knock on the door.

← Оригинальное предложение (source)

RU: В дверь постучали.

← Перевод (target)

DE: Ein Klopfen an der Tür **unterbrach** ihn.

ITA: **Fu interrotto** da qualcuno che bussava alla porta.

ES: **Le interrumpieron** unos golpes en la puerta.

CZ: **Přerušilo ho** zaklepaní na dveře.

BG: **Прекъсна го** почукване на вратата.

SE: **Han avbröts** av en knackning på dörren.

FR: **Il fut interrompu** par des coups frappés à la porte.

Переводные единицы

Предложение (или несколько предложений) на языке оригинала с соответствующими ему переводами на представленные в корпусе языки называется **переводной единицей**.

EN: **He was interrupted** by a knock on the door.

RU: В дверь постучали.

DE: Ein Klopfen an der Tür **unterbrach** ihn.

ITA: **Fu interrotto** da qualcuno che bussava alla porta.

ES: **Le interrumpieron** unos golpes en la puerta.

CZ: **Přerušilo ho** zaklepaní na dveře.

BG: **Прекъсна го** почукване на вратата.

SE: **Han avbröts** av en knackning på dörren.

FR: **Il fut interrompu** par des coups frappés à la porte.

← Оригинальное предложение (source)

← Перевод (target)

Переводная единица

Переводные эквиваленты

Языковые выражения, используемые при переводе, которые соответствуют словам и конструкциям из оригинального текста, называются **переводными эквивалентами**.

EN: He was interrupted by a knock on the door.

RU: В дверь постучали.

DE: Ein Klopfen an der Tür unterbrach ihn.

ITA: Fu interrotto da qualcuno che bussava alla porta.

ES: Le interrumpieron unos golpes en la puerta.

CZ: Přerušilo ho zklepání na dveře.

BG: Прекъсна го почукване на вратата.

SE: Han avbröts av en knackning på dörren.

FR: Il fut interrompu par des coups frappés à la porte.

Переводные эквиваленты

Наполнение корпуса

Важным атрибутом параллельных корпусов является наличие в них соответствий различных частей корпуса друг другу.

Идеальный параллельный корпус выглядит примерно так:

	RU	EN	FR	IT	DE
RU	1.5M	1.5M	1.5M	1.5M	1.5M
EN	1.5M	1.5M	1.5M	1.5M	1.5M
FR	1.5M	1.5M	1.5M	1.5M	1.5M
IT	1.5M	1.5M	1.5M	1.5M	1.5M
DE	1.5M	1.5M	1.5M	1.5M	1.5M

Наполнение корпуса

Важным атрибутом параллельных корпусов является наличие в них соответствий различных частей корпуса друг другу.

Идеальный параллельный корпус выглядит примерно так:

	RU	EN	FR	IT	DE
RU	1.5M	1.5M	1.5M	1.5M	1.5M
EN	1.5M	1.5M	1.5M	1.5M	1.5M
FR	1.5M	1.5M	1.5M	1.5M	1.5M
IT	1.5M	1.5M	1.5M	1.5M	1.5M
DE	1.5M	1.5M	1.5M	1.5M	1.5M

Например, взяли несколько произведений и для каждого нашли переводы для всех языков, представленных в корпусе

[illegible]

- Не существует каких-то переводов, либо они недоступны
- Очень специфические тексты становятся массово параллельными

Как обстоят дела с параллельными корпусами?

- С одной стороны ресурсов, предлагающих исследователям в открытом доступе мультязычные параллельные тексты имеется немало.
- Однако по факту напрямую использовать в исследованиях можно не все из них, это связано с
 - неоднородностью объема соответствий переводных пар в рамках одного корпуса,
 - отсутствием унифицированной разметки,
 - спецификой самих текстов, входящих в корпус
- Двужычных параллельных корпусов существует гораздо больше. В принципе, чем меньше в корпусе языков, тем больше шансов найти корпус большого объема.

Исследования на материале параллельных корпусов

Вариативность переводных эквивалентов

- Прелесть работы с переводами состоит в том, что изучая некоторую конструкцию или даже просто употребление слова можно найти много нетривиальных соответствий.
- Соответствия не всегда однородные, не всегда прилагательному соответствует прилагательное, а существительному существительное, пассиву соответствует пассив или актив.

Примеры

Иногда это обусловлено наличием в языках нескольких “официально” зафиксированных вариантов выражения, например, для каузативов:

– Lexical causatives, e. g., *break_{TR}* or *walk_{TR}*;

Примеры

Иногда это обусловлено наличием в языках нескольких “официально” зафиксированных вариантов выражения, например, для каузативов:

- Lexical causatives, e. g., *break_{TR}* or *walk_{TR}*;
- Morphological causatives, e. g., *tone change*, *reduplication*, or *affixation*;

Примеры

Иногда это обусловлено наличием в языках нескольких “официально” зафиксированных вариантов выражения, например, для каузативов:

- Lexical causatives, e. g., *break_{TR}* or *walk_{TR}*;
- Morphological causatives, e. g., *tone change*, *reduplication*, or *affixation*;
- Complex predicates, e. g., serial verbs, French *faire* ‘make’ + VINFINF, or causative particles;

Примеры

Иногда это обусловлено наличием в языках нескольких “официально” зафиксированных вариантов выражения, например, для каузативов:

- Lexical causatives, e. g., *break_{TR}* or *walk_{TR}*;
- Morphological causatives, e. g., *tone change*, *reduplication*, or *affixation*;
- Complex predicates, e. g., serial verbs, French *faire* ‘make’ + VINF, or causative particles;
- Periphrastic causatives, where the causatives are represented by verbs that belong to separate clauses, e. g., French *laisser* ‘let’ + NP + VINF or Portuguese *fazer* ‘make’ + (NP) + VINF.

Примеры

Пассивные конструкции

- Соответствие лексической единице

- EN: **Her hands shook** slightly

ITA: **Le sue mani erano scosse da un lieve tremito**

- Соответствие конструкции вроде *there is/are*

- EN: **There was a dull murmur** of assent **throughout the class**.

ITA: **La classe fu percorsa da un cupo mormorio** di assenso.

- EN: **There were movements** from the watching crowd **in front of the castle**, <...>

ITA: **La folla davanti al castello fu percorsa da un fremito** <...>

Персоны

Персоны

Bernhard Wälchli, Michael Cysouw

Wälchli B., Cysouw M. Lexical typology through similarity semantics: Toward a semantic map of motion verbs // *Linguistics*. 2012. No 3 (50). C. 671–710.

Östen Dahl

Dahl Ö. From questionnaires to parallel corpora in typology // *Language Typology and Universals*. 2007. No 2 (60). C. 172–181



Michael Cysouw



Östen Dahl

Персоны

Bernhard Wälchli, Michael Cysouw

Wälchli B., Cysouw M. Lexical typology through similarity semantics: Toward a semantic map of motion verbs // *Linguistics*. 2012. No 3 (50). C. 671–710.

Östen Dahl

Dahl Ö. From questionnaires to parallel corpora in typology // *Language Typology and Universals*. 2007. No 2 (60). C. 172–181



Michael Cysouw



Östen Dahl

Персоны: Bernhard Wälchli, Michael Cysouw, Östen Dahl

Одними из первых указали на преимущества параллельных корпусов

Персоны: Bernhard Wälchli, Michael Cysouw, Östen Dahl

Одними из первых указали на преимущества параллельных корпусов, а именно:

- Охват сразу с нескольких языков
- Возможность применения различных современных методов анализа.
- Работа с ситуациями в одинаковом контекстном окружении. То есть исследователь располагает примерами с одинаковой семантической и прагматической составляющей. Это не разрозненные примеры из разных источников, это сопоставимый материал.

Персоны: Bernhard Wälchli, Mi

Одними из первых указали на преим

- Охват сразу с нескольких языков
- Возможность применения различны

Volume 60 Issue 2, Issue of STUF - Language Typology and Universals

Это не разрозненные примеры из р:

Requires Authentication September 25, 2009

Parallel texts: using translational equivalents in linguistic typology

Michael Cysouw, Bernhard Wälchli

Page range: 95-99

More ▼

Cite this

Requires Authentication September 25, 2009

Harry Potter meets *Le petit prince* – On the usefulness of parallel corpora in crosslinguistic investigations

Thomas Stolz

Page range: 100-117

More ▼

Cite this

Requires Authentication September 25, 2009

Advantages and disadvantages of using parallel texts in typological investigations

Bernhard Wälchli

Page range: 118-127

More ▼

Cite this

Page range: 135-147

More ▼

Cite this

Requires Authentication September 25, 2009

Some remarks on the use of Bible translations as parallel texts in linguistic research

Lourens de Vries

Page range: 148-157

More ▼

Cite this

Requires Authentication September 25, 2009

Using Strong's Numbers in the Bible to test an automatic alignment of parallel texts

Michael Cysouw, Chris Biemann, Matthias Ongyerth

Page range: 158-171

More ▼

Cite this

ей.

Персоны



- **Robert Östling, University of Stockholm**

- <https://www.su.se/english/profiles/robe-1.187515>
- Важно не перепутать с экономистом-тёзкой!
- Östling R. 6. Studying colexification through massively parallel corpora Berlin, Boston: De Gruyter, 2016
- Östling R. [Word order typology through multilingual word alignment](#), In: The 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing: Proceedings of the Conference, Volume 2: Short Papers, 2015, p. 205-211

- **Jörg Tiedemann, University of Helsinki**

- <https://researchportal.helsinki.fi/en/persons/j%C3%B6rg-tiedemann>
- Отец OPUS'a
- Östling R., Tiedemann J. Efficient Word Alignment with Markov Chain Monte Carlo // The Prague Bulletin of Mathematical Linguistics. 2016. No (106). C. 125–146



Персоны

- **Ruprecht von Waldenfels, University of Jena**

- <https://www.gw.uni-jena.de/fakultaet/institut-fuer-slavistik-und-kaukasu/sstudien/mitarbeiterinnen/von-waldenfels-ruprecht>
- Славист, создатель корпуса ParaSOL
- Waldenfels R. von Compiling a parallel corpus of Slavic languages. Text strategies, tools and the question of lemmatization in alignment // Beiträge der europäischen slavistischen Linguistik (POLYSLAV). 2006. (9). С. 123–138.



- **Дмитрий В. Сичинава, ИРЯ РАН**

- <https://www.ruslang.ru/publica/sichinava>
- Типолог, специалист по славянским языкам, один из создателей параллельных корпусов в составе НКРЯ
- Sitchinava D. V, others Parallel corpora within the Russian National Corpus // Prace Filologiczne. 2012. No 63. С. 271–278



Персоны

- **Наталия Лёвшина, MPI for Psycholinguistics Nijmegen**

- Создательница корпуса ParTY, развивает применение количественных методов и параллельных корпусов в типологии
- Levshina N. Why we need a token-based typology: A case study of analytic and lexical causatives in fifteen European languages // Folia Linguistica. 2016a. No 2 (50).



- **Сергей С. Сай, ИЛИ РАН**

- Типолог, специалист по аргументной структуре
- Say, Sergey. [Nominal causal constructions across Slavic: semantic contrasts in a parallel corpus perspective.](#) Slavia, 2021, 90, 2. P. 182–201.



Методы анализа

- **Построение семантических карт.** Многомерное шкалирование позволяет сгруппировать переводные единицы, с одинаковыми свойствами на основе их сходств в маркировании.
- **Сетевой анализ** дает возможность проанализировать, как взаимодействуют различные способы выражения, без привязки к переводным единицам.
- При помощи **факторного анализа** можно оценивать влияние различных признаков на выбор употребления конструкции.

Перенос разметки

- Наличие выравнивания в параллельном корпусе позволяет осуществлять перенос разметки и использовать это в исследованиях.
 - Например, можно разметить семантические признаки или характеристики ситуации (одушевленность участников, наличие негативного воздействия и др.) и перенести эту информацию с одного языка на другой.

Перенос разметки

- Наличие выравнивания в параллельном корпусе позволяет осуществлять перенос разметки и использовать это в исследованиях.
 - Например, можно разметить семантические признаки или характеристики ситуации (одушевленность участников, наличие негативного воздействия и др.) и перенести эту информацию с одного языка на другой.
- Перенос разметки позволяет упростить процесс аннотирования мультязычного корпуса: достаточно разметить лишь часть данных.

Перенос разметки

- Наличие выравнивания в параллельном корпусе позволяет осуществлять перенос разметки и использовать это в исследованиях.
 - Например, можно разметить семантические признаки или характеристики ситуации (одушевленность участников, наличие негативного воздействия и др.) и перенести эту информацию с одного языка на другой.
- Перенос разметки позволяет упростить процесс аннотирования мультязычного корпуса: достаточно разметить лишь часть данных.
- Плохо сработает для грамматических характеристик (иначе нам не нужно было бы использовать кучу парсеров)

S. Say, Nominal causal constructions across Slavic: semantic contrasts in a parallel corpus perspective

С. Сай в своей работе [Say 2021b] рассматривает именные причинные конструкции, то есть такие причинные конструкции, в которых событие-причина выражена именной группой.

(1) Маргарита ... заплакала [**от** бол-**и** в руке и ноге]. (Russian)

‘Margarita ... burst into tears from the pain in her arm and leg.’

(2) Я погибаю [**из-за** любв-**и**]. (Russian)

‘I’m perishing on account of love.’

Автор анализирует способы маркирования причинных именных групп в нескольких славянских языках. При помощи семантической карты получилось рассмотреть, как в рамках единого семантического пространства распределены предложные показатели, используемые в языках выборки.

S. Say, Nominal causal constructions across Slavic: semantic contrasts in a parallel corpus perspective

С. Сай в своей работе [Say 2021b] рассматривает такие причинные конструкции, в частности, в к

- (1) Маргарита ... заплакала [от бол-
'Margarita ... burst into tears from the
(2) Я погибаю [из-за любв-и]. (Russ
'I'm perishing on account of love.

Автор анализирует способы маркировки причинности в славянских языках. При помощи семантических рамок единого семантического пространства он анализирует используемые в языках выборки.

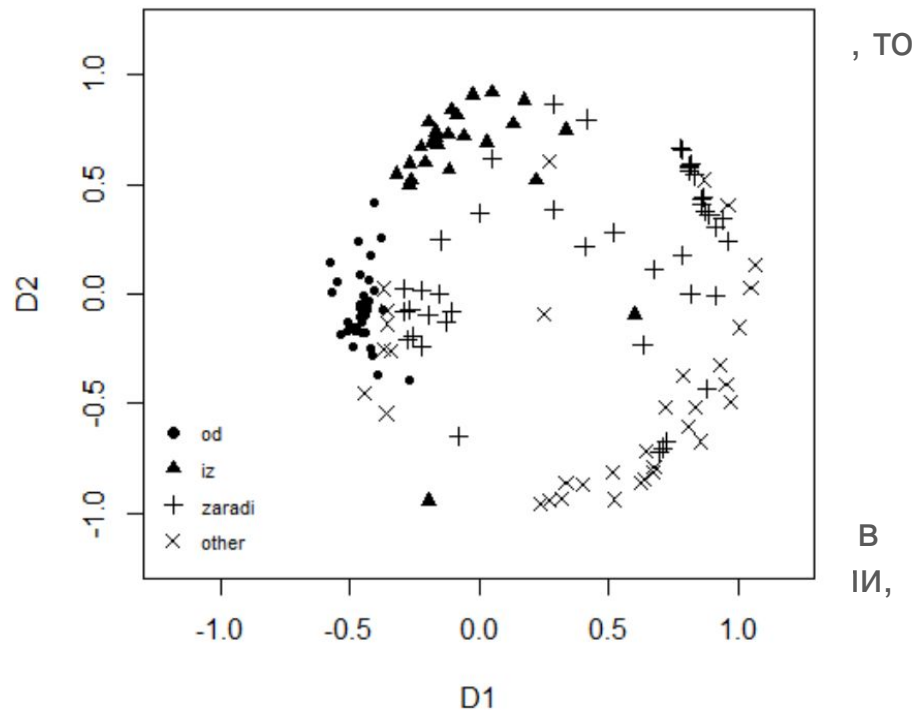


Fig. 7. Bulgarian nominal causal markers

Natalia Levshina, Why we need a token-based typology: A case study of analytic and lexical causatives in fifteen European languages

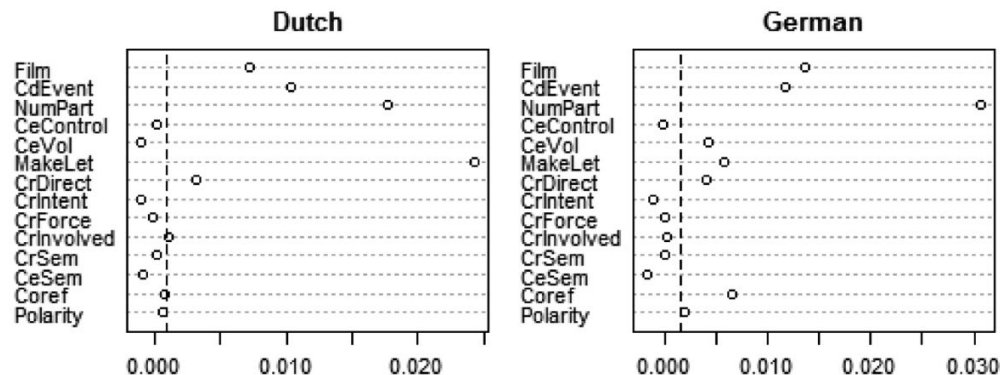
Table 3: Parameters of variation of lexical and analytic causatives operationalized as variables.

	Variable	Abbreviation	Values	Expectations
1	Aktionsart of the caused event	<i>CdEvent</i>	“Nonaction” “Action”	Lexical Analytic
2	Number of main participants	<i>NumPart</i>	“2” “3”	Lexical Analytic
3	Control of Causee	<i>CeControl</i>	“Yes” “No”	Analytic Lexical
4	Causee acting willingly	<i>CeVol</i>	“Yes” “No” “Undef”	Analytic Lexical No clear expectations
5	Making or letting	<i>MakeLet</i>	“Make” “Let”	Lexical Analytic
6	Causer acting directly	<i>CrDirect</i>	“Yes” “No”	Lexical Analytic
7	Causer acting intentionally	<i>CrIntent</i>	“Yes” “No”	Lexical Analytic
8	Causer acting forcefully	<i>CrForce</i>	“Yes” “No”	Analytic Lexical
9	Causer involved in caused event	<i>CrInvolved</i>	“Yes” “No”	No clear expectations
10	Semantics of Causer	<i>CrSem</i>	“Anim” “Inanim”	Lexical Analytic
11	Semantics of Causee	<i>CeSem</i>	“Anim” “Inanim”	Analytic Lexical
12	Coreferentiality of Causer with other main participants	<i>Coref</i>	“Yes” “No”	No clear expectations
13	Polarity	<i>Polarity</i>	“Pos” “Neg”	No clear expectations

Natalia Levshina, Why we need a token-based typology: A case study of analytic and lexical causatives in fifteen European languages

Table 3: Parameters of variation of lexical and analytic causatives operationalized as variables.

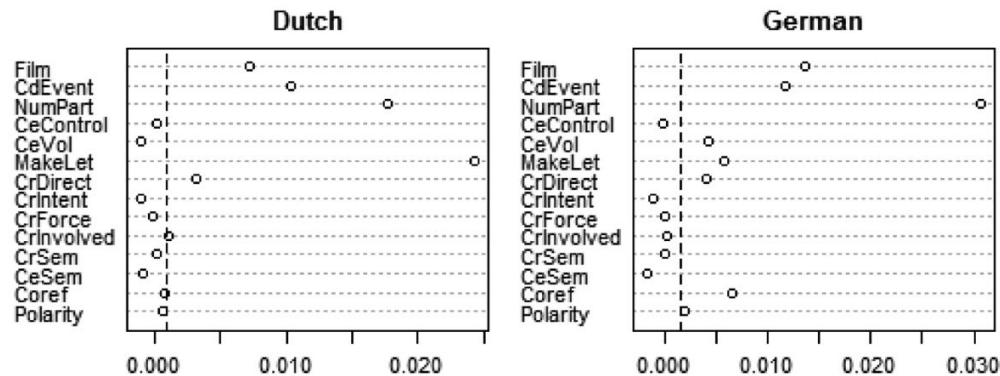
	Variable	Abbreviation	Values	Expectations
1	Aktionsart of the caused event	<i>CdEvent</i>	“Nonaction” “Action”	Lexical Analytic
2	Number of main participants	<i>NumPart</i>	“2” “3”	Lexical Analytic
3	Control of Causee	<i>CeControl</i>	“Yes” “No”	Analytic Lexical
4	Causee acting willingly	<i>CeVol</i>	“Yes” “No” “Undef”	Analytic Lexical No clear expectations
5	Making or letting	<i>MakeLet</i>	“Make” “Let”	Lexical Analytic
6	Causer acting directly	<i>CrDirect</i>	“Yes” “No”	Lexical Analytic
7	Causer acting intentionally	<i>CrIntent</i>	“Yes” “No”	Lexical Analytic
8	Causer acting forcefully	<i>CrForce</i>	“Yes” “No”	Analytic Lexical
9	Causer involved in caused event	<i>CrInvolved</i>	“Yes” “No”	No clear expectations
10	Semantics of Causer	<i>CrSem</i>	“Anim” “Inanim”	Lexical Analytic
11	Semantics of Causee	<i>CeSem</i>	“Anim” “Inanim”	Analytic Lexical
12	Coreferentiality of Causer with other main participants	<i>Coref</i>	“Yes” “No”	No clear expectations
13	Polarity	<i>Polarity</i>	“Pos” “Neg”	No clear expectations



Natalia Levshina, Why we need a token-based typology: A case study of analytic and lexical causatives in fifteen European languages

Table 3: Parameters of variation of lexical and analytic causatives operationalized as variables.

	Variable	Abbreviation	Values	Expectations
1	Aktionsart of the caused event	<i>CdEvent</i>	“Nonaction” “Action”	Lexical Analytic
2	Number of main participants	<i>NumPart</i>	“2” “3”	Lexical Analytic
3	Control of Causee	<i>CeControl</i>	“Yes” “No”	Analytic Lexical
4	Causee acting willingly	<i>CeVol</i>	“Yes” “No” “Undef”	Analytic Lexical No clear expectations
5	Making or letting	<i>MakeLet</i>	“Make” “Let”	Lexical Analytic
6	Causer acting directly	<i>CrDirect</i>	“Yes” “No”	Lexical Analytic
7	Causer acting intentionally	<i>CrIntent</i>	“Yes” “No”	Lexical Analytic
8	Causer acting forcefully	<i>CrForce</i>	“Yes” “No”	Analytic Lexical
9	Causer involved in caused event	<i>CrInvolved</i>	“Yes” “No”	No clear expectations
10	Semantics of Causer	<i>CrSem</i>	“Anim” “Inanim”	Lexical Analytic
11	Semantics of Causee	<i>CeSem</i>	“Anim” “Inanim”	Analytic Lexical
12	Coreferentiality of Causer with other main participants	<i>Coref</i>	“Yes” “No”	No clear expectations
13	Polarity	<i>Polarity</i>	“Pos” “Neg”	No clear expectations



Подход:

- 1) Размечаем признаки на материале одного языка
- 2) Переносим это на остальные языки
- 3) Строим модели для языков выборки на основе одинаковых признаков, но разных способов маркирования ситуаций

Создание корпусов

Применение

Преимущества параллельных корпусов

- Охват сразу с нескольких языков
- Возможность применения различных современных методов анализа.
- Работа с ситуациями в одинаковом контекстном окружении. То есть исследователь располагает примерами с одинаковой семантической и прагматической составляющей. Это не разрозненные примеры из разных источников, это сопоставимый материал.

Применение

Преимущества параллельных корпусов

- Охват сразу с нескольких языков
- Возможность применения различных современных методов анализа.
- Работа с ситуациями в одинаковом контекстном окружении. То есть исследователь располагает примерами с одинаковой семантической и прагматической составляющей. Это не разрозненные примеры из разных источников, это сопоставимый материал.

Применение:

- Сравнительные исследования, типология
- Машинный перевод

Чего мы хотим от параллельных корпусов?

=Что придется реализовать технически

Чего мы хотим от параллельных корпусов?

=Что придется реализовать технически

- Много разных языков
- Много текстов и соответствий
- Несколько этапов обработки

Pipeline создания параллельного корпуса

Сбор текстов (проверка, что все со всем параллельно)

Предобработка текстов (разбиение на предложения, токенизация)

Выравнивание по предложениям

Выравнивание по словам

Разметка

Разметка и выравнивание

Какие проблемы могут возникнуть?

Как можно решать эти проблемы?

Разметка и выравнивание

Какие проблемы могут возникнуть?

Как можно решать эти проблемы?

- Нет универсальных инструментов, а нам нужно обработать очень разнообразные данные

Важные понятия

Выравнивание (alignment)

Переводная единица (translational unit)

Переводные эквиваленты (translational equivalents)

Язык оригинала (source language)

Язык перевода (target language)