

Выравнивание текстов в параллельном корпусе: подходы и инструменты

Рекап

1. Определение параллельного корпуса
2. Переводные единицы

Рекап

1. Определение параллельного корпуса
2. Переводные единицы

eng	And here is the joke. A day earlier, she'd called him to ask for my phone number. [André Aciman. Find Me (2019) A
	снята] ←...→
rus	И вот что забавно: днем ранее она связалась с ним, чтобы узнать мой номер. [André Aciman. Find Me (2019)
	снята] ←...→

3. Переводные эквиваленты

Рекап

1. Определение параллельного корпуса
2. Переводные единицы

eng	And here is the joke. A day earlier, she'd called him to ask for my phone number. [André Aciman. Find Me (2019) A
	снята] ←...→
rus	И вот что забавно: днем ранее она связалась с ним, чтобы узнать мой номер. [André Aciman. Find Me (2019)
	снята] ←...→

3. Переводные эквиваленты и нетривиальные соответствия в переводе
4. Pipeline создания параллельного корпуса

Примеры параллельных текстов

Примеры параллельных текстов

Наиболее известные:

- Библия
- Гарри Поттер
- Маленький принц
- Винни Пух

Чуть менее тривиальные:

- Инструкции
- Субтитры к фильмам

Примеры параллельных корпусов

- Параллельный корпус для Слова о полку Игореве
- Выровненный программными средствами текст октоиха

Ἦχος α.
ΤΩ ΣΑΒΒΑΤΩ ΕΣΠΕΡΑΣ, ΕΝ ΤΩ ΜΙΚΡΩ Ε-
ΣΠΕΡΙΝΩ

Στίχ. Ἀπὸ φυλακῆς πρωίας μέχρι νυκτός, ἀπὸ
φυλακῆς πρωίας, ἐλπισάτω Ἰσραὴλ ἐπὶ τὸν Κύριον.

Τὰς ἐσπερινὰς ἡμῶν εὐχὰς, πρόσδεξαι Ἁγίε Κύ-
ριε, καὶ παράσχου ἡμῖν ἄφεσιν ἁμαρτιῶν, ὅτι μό-
νος εἶ ὁ δείξας ἐν κόσμῳ τὴν Ἀνάστασιν. (Δίς).

Κυκλώσατε λαοὶ Σιών, καὶ περιλάβετε αὐτήν,
καὶ δότε δόξαν ἐν αὐτῇ, τῷ Ἀναστάντι ἐκ νεκρῶν,
ὅτι αὐτός ἐστιν ὁ Θεὸς ἡμῶν, ὁ λυτρωσάμενος
ἡμᾶς ἐκ τῶν ἀνομιῶν ἡμῶν.

Кни́га моле́бна
съ бѣгомъ стѣмъ ѡсмогласника, содержащаа
въ себѣ подобающее возслѣдованіе воскрѣныхъ
слѣжбы ѡсмй гласѡвъ съ шестію днѣй

ГЛАСЬ А

Въ съббоѡтъ вечера, на малѣй вечерни,

на Гди воззвахъ, поста́вимъ стѣхѡвъ дѣ: и по-
емъ стѣхиры воскрѣны ѡсмогласника г, повтор-
ающе а-й стѣхъ, гласъ а. Твореніе прѣбнаго о́тца
нашегого іωάνна дамаскіна

Стѣхъ: Ѡ стражи ѹтренніа до нощи, ѡ стра-
жи ѹтренніа, да ѹпова́еть іиѡль на гда.

Вечерніа наша мѣтвы, прѣими стѣй гди, и
подаждь намъ ѡставленіе грѣхѡвъ, іако ѡдинъ
еси іавлѣй въ мірѣ воскрѣне.

Ѡбыдите людіе сіѡнь, и ѡбимите єго, и да-
дите славу въ немъ воскрѣшемъ изъ мѣртвыхъ:
іако то́й єсть бгъ нашъ, избавлѣй насъ ѡ без-
законіи нашихъ.

Примеры корпусов

1. Коллекция [OPUS](#) - платформа для работы с большим количеством языков, которая предлагает множество различных инструментов
2. [ParTY](#) (субтитры к фильмам)
3. Параллельный корпус в составе НКРЯ (художественные произведения)
4. [Библия](#)

Азбука веры » Библия » Книга Бытия » гл.1



Книга Бытия, глава 1 ☆



Толкования главы

Синодальный : × 🔊

Греческий : ×

Языки +

1:1 В начале сотворил Бог небо и землю.

1:2 Земля же была безвидна и пуста, и тьма над бездною, и Дух Божий носился над водою.

1:3 И сказал Бог: да будет свет. И стал свет.

1:4 И увидел Бог свет, что он хорош, и отделил Бог свет от тьмы.

<input type="checkbox"/> ἐν ἀρχῇ ἐποίησεν ὁ θεὸς τὸν οὐρανὸν καὶ τὴν γῆν	<input type="checkbox"/>
<input type="checkbox"/> ἡ δὲ γῆ ἦν ἀόρατος καὶ ἀκατασκεύαστος καὶ σκότος ἐπάνω τῆς ἀβύσσου καὶ πνεῦμα θεοῦ ἐπεφέρετο ἐπάνω τοῦ ὕδατος	<input type="checkbox"/>
<input type="checkbox"/> καὶ εἶπεν ὁ θεός γενεθήτω φῶς καὶ ἐγένετο φῶς	<input type="checkbox"/>
<input type="checkbox"/> καὶ εἶδεν ὁ θεὸς τὸ φῶς ὅτι καλόν καὶ διεχώρισεν ὁ θεὸς ἀνὰ μέσον τοῦ φωτὸς καὶ ἀνὰ μέσον τοῦ σκότους	<input type="checkbox"/>

Примеры корпусов

Корпус слушаний Европарламента

- 21 официальный язык ЕС
- Все подкорпуса выровнены по английскому
- XML, размечены говорящие, файл соответствует дню слушаний
- Свободный для скачивания

Корпус европейского права (The JRC-Acquis Multilingual Parallel Corpus)

- Действующее право ЕС. 22 языка
- Общий объём – 1 млрд слов
- Автоматическое выравнивание (венгерская программа HunAlign).

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

Выравнивание

Установление соответствия между единицами (абзацем, предложениями, словами) оригинального и переводного текста.

Основные подходы	Проблемы
→ пословное выравнивание	несовпадение лексем в разных языках
→ выравнивание по предложениям	количество предложений и их расположение в тексте может не совпадать (перекрестная структура текста)

Выравнивание по предложениях vs По словам

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 эквиваленты

Выравнивание по словам - основной принцип

Предложение языка L1: $s(1..j) = s_1 + \dots + s_j$ $s \rightarrow \text{source}$

Предложение языка L2: $t(1..i) = t_1 + \dots + t_i$ $t \rightarrow \text{target}$

Под выравниванием понимается отображение множества позиций слов исходного предложения $\{1 \dots j\}$ во множество позиций слов целевого $\{1 \dots i\}$ предложения:

$$a: i \rightarrow j$$

Выравнивание по словам

Пример пословного выравнивания с помощью программы Giza++:

Ls: Штифт (1) для (2) использования (3) в (4) стоматологии (5).

Lt: NULL(0) Implant(1) à (2) usage(3) dentaire(4).

a: { 1 → 1, 2 → 0, 3 → 3, 4 → 0, 4 → 5 }

	1	2	3	4	5
1					
2					
3					
4					

Потренируемся :)

Ls: Он(1) будет(2) рад(3) цветам(4).

Lt: NULL(0) He(1) will(2) love(3) the(4) flowers(5).

Ls: Луна(1) освещала(2) голые(3) обугленные(4) стволы(5).

Lt: NULL(0) The(1) moon(2) illuminated(3) bare(4), charred(5) tree(6) trunks(7).

Ls: У(1) меня(2) красивые(3) руки(4).

Lt: NULL(0) I(1) have(2) nice(3) hands(4).

Основные подходы к пословному выравниванию

1. Эвристические модели: для каждой пары соответствующих друг другу предложений строится матрица со значением ассоциативных мер между каждым словом исходного предложения и каждым словом целевого предложения. Для выравнивания со словом исходного предложения выбирается слово с большей мерой.
2. Самообучающиеся модели: модели, основанные на принципах вероятности (в их основе - машинное обучение и скрытые марковские цепи)

Выравнивание по словам - литература:

- Östling R., Tiedemann J. [Efficient Word Alignment with Markov Chain Monte Carlo](#) // The Prague Bulletin of Mathematical Linguistics. 2016. No 1 (106). С. 125–146
- 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

Выравнивание по предложениям

В наиболее простых случаях для деления текста на предложения используются элементарные синтаксические признаки конца и начала предложения – пунктуация, использование заглавных букв, знаки абзаца.

Основные подходы:

1. По длине ([LF-aligner](#), Евклид - оболочки на базе [HunAlign](#))
2. По лексике (алгоритмы дистрибутивной семантики, BERT, Lingtrain)

HunAlign. Выравнивание по предложениям

In the absence of a dictionary, it first falls back to sentence-length information, and then builds an automatic dictionary based on this alignment. Then it realigns the text in a second pass, using the automatic dictionary.

Like most sentence aligners, hunalign does not deal with changes of sentence order: it is unable to come up with crossing alignments, i.e., segments A and B in one language corresponding to segments B' A' in the other language.



В HunAlign подгружаются лемматизированные файлы

Выравнивание по предложениям

«Евклид» (оболочка для HunAlign). Тексты выравниваются попарно и затем «склеиваются» в единый XML в соответствии с разделением предложений в оригинале:

```
<para id="2">
```

```
  <se lang="ru" variant_id="0">Марта 25 числа случилось в Петербурге  
необыкновенно странное происшествие.</se>
```

```
  <se lang="fr" variant_id="1">Le 25 mars, un événement tout à fait étrange  
s'est produit à Pétersbourg.</se>
```

```
  <se lang="fr" variant_id="2">Ce jour-là, 25 mars dernier, Pétersbourg fut le  
théâtre d'une aventure des plus étranges.</se>
```

```
  <se lang="fr" variant_id="3">Le 25 mars il est arrivé à Pétersbourg un  
événement extrêmement bizarre.</se>
```

```
</para>
```

Выравнивание по предложениям. LF-aligner

74 PTK

Specify the languages of your texts:

Number of languages (usually, 2): 2

Language 1: English

Language 2: Hungarian

Note: you can change the default languages by editing LF_aligner_

1	въ срѣда вечерѣ	ὁ Περμπτη ὁ ς ' Εβδομαδος	lemmas_or_rnc-lemmas_greek_perseus
2	празднити вознесение господь Богъ и спасъ нашъ Иисусъ Христъ	ὁ ὁ Μικρῶι Εσπερις Στιχηρὸς ὁ ἑορτῇ	lemmas_or_rnc-lemmas_greek_perseus
3	часть 9-и трипсалмныи	ἤχος πῖ β	lemmas_or_rnc-lemmas_greek_perseus
4		ὁ Κύριος ἀλαμβάνω εἰς οὐρανός , ἵνα πέμψω ὁ Παράκλητος ὁ κόσμος , ὁ οὐρανός ἐτοιμάζω ὁ θρόνος αὐτός , νεφέλα ὁ ἐπίβασις αὐτός , Ἀγγελοι θαυμάζω , ἀνθρωπος ὁράω ὑπεράνω αὐτός , ὁ Πατήρ ἐκδέχομαι , ὅς ἐν κόλπος ἐχω συναΐδιος	lemmas_or_rnc-lemmas_greek_perseus
5	по начало глагоити чтець трисвятой , и до недѣля 8-я во всякой церковный пения на малый вечерня	ὁ Πνεῦμα ὁ ἅγιος κελεῦω πᾶς ὁ Ἀγγέλοι αὐτός . Ἀρατε πύλη ὁ ἀρχων ἐγώ , πᾶς ὁ ἐθνος κροτάω χεῖρ	lemmas_or_rnc-lemmas_greek_perseus
6	на господь воззвати : поставити стихирь 4 , гласъ 6	ὅτι ἀνέβη Χριστός , ὅπου εἰμί ὁ πρότερος ἤχος δέ	lemmas_or_rnc-lemmas_greek_perseus
7	Господь вознестися на небо : господь твой вознесение : на гора святой : господь , смотрѣние совершивъ таинство : зърѣти ниже на великий вечерни	Κύριος , ὁ σὸς , Ἀναλήψης , Εξεπλάγησας ὁ Χερουβίμ , θεωρήζω σύ ὁ θεόν , ἐπὶ νεφέλῃ ἀνέρχομαι , ὁ εἰ ὅς αὐτός καθέζομαι , καὶ δοξάζω σύ , ὅτι χρηστός ὁ ἐλεός σου , δόξα σύ ἤχος πῖ	lemmas_or_rnc-lemmas_greek_perseus
8	слава , и нынѣ : гласъ тыже : господь , апостоль яко видѣти тя	β ἐν ὁ ὁρος ὁ ἅγιός , θεωρέω σύ ὁ ὑψὼ Χριστός , ὁ ἀπαύγασμα ὁ δόξα ὁ Πατήρ , ἀνυμνέω σύ ὁ φωτειδής ὁ πρόσωπον μορφῇ , προσκυνέω σύ ὁ πάθημα , τιμάω	lemmas_or_rnc-lemmas_greek_perseus

Merge (F1)Split (F2)Shift up (F3)Shift down (F4)

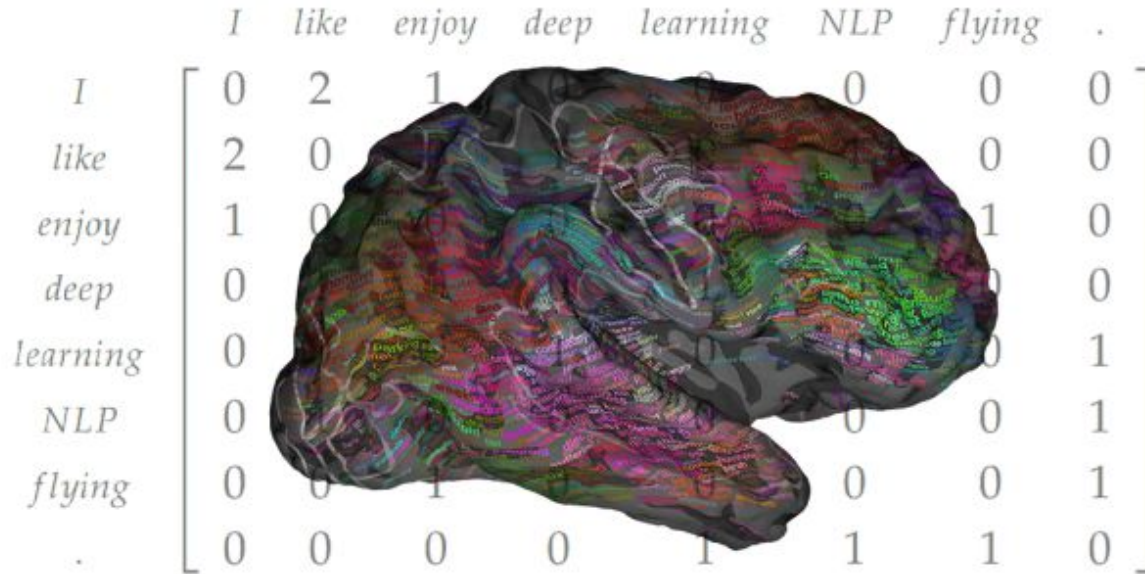
Дистрибутивная семантика

Дистрибутивная семантика - это метод исследования языка, основанный на дистрибуции отдельных единиц в тексте и не использующий сведений о конкретном лексическом и грамматическом значении слов.

Впервые предложен Л.Блумфильдом в 20ые гг. XX века и применялся, главным образом, в фонологии и морфологии. Используя контексты в качестве исходных данных (какие единицы могут взаимодействовать друг с другом), можно выделить основные единицы языка, объединить в классы и установить отношения встречаемости между ними.

Дистрибутивная семантика позволяет автоматически устанавливать контексты употребления слова и устанавливать семантические связи между лексемами.

Дистрибутивная семантика



We want a **machine** to imitate human brain and **understand meaning of words**.

Дистрибутивная семантика

В качестве инструмента дистрибутивной семантики используется линейная алгебра. Информация о дистрибуции единиц представляется в виде многомерных векторов.

Сами векторы соответствуют лексическим единицам (словам или словосочетаниям), а измерения векторов - контекстам, в которых эти единицы встречаются.

$$A_{m,n} = \begin{matrix} & & w_1 & \text{drink} & w_3 & \dots & w_m \\ \text{coffee} & & 0 & 1 & 0 & \dots & 1 \\ w_2 & & 1 & 0 & 2 & \dots & 0 \\ \text{tea} & & 0 & 2 & 0 & \dots & 3 \\ \vdots & & \vdots & \vdots & \vdots & \ddots & \vdots \\ w_m & & 0 & 0 & 1 & \dots & 0 \end{matrix}$$

Словесное пространство для слов английского языка *coffee* и *tea*

Дистрибутивная семантика

Семантическая близость между единицами вычисляется как косинусное расстояние между векторами.

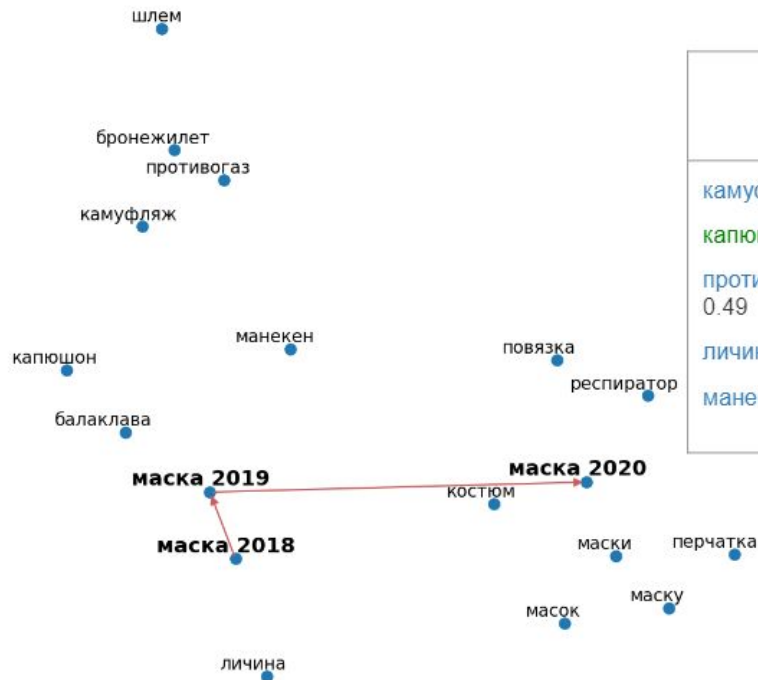
$$\cos(w1, w2) = \frac{\vec{V}(w1) \times \vec{V}(w2)}{|\vec{V}(w1)| \times |\vec{V}(w2)|}$$

$$\cos(\text{tomat}, \text{philosophy}) = 0.00698$$

$$\cos(\text{pomidor}, \text{philosophy}) = -0.03429$$

$$\cos(\text{tomat}, \text{pomidor}) = 0.65049$$

Дистрибутивная семантика



2018	2019	2020
		Есть сдвиг! (0.64)
камуфляж NOUN 0.52	камуфляж NOUN 0.57	повязка NOUN 0.60
капюшон NOUN 0.51	капюшон NOUN 0.50	одноразовый ADJ 0.55
противогаз NOUN 0.49	перчатка NOUN 0.47	костюм NOUN 0.50
личина NOUN 0.47	шлем NOUN 0.45	защитный ADJ 0.49
манекен NOUN 0.45	противогаз NOUN 0.45	перчатка NOUN 0.46

<https://shiftry.rusvectors.org/en/about/>

Дистрибутивная семантика

С помощью дистрибутивной семантики легко и быстро решается множество проблем: снятие семантической неоднозначности, тематическая кластеризация, генерация тезаурусов, поиск синонимов, антонимов, гиперонимов и многое другое.

Чем полезно для выравнивания по предложениям?

Дистрибутивная семантика и выравнивание текстов

Lingtrain Aligner: <https://github.com/averkij/lingtrain-aligner>

В чем идея:

→ Модель берет заданное количество строк первого текста и подбирает в соответствующем фрагменте второго текста лучшие соответствия, используя векторные представления.

[Ὁ οἶκος]

[Τὸν πρὸ ἡλίου Ἥλιον, δύναντα ποτὲ ἐν τάφῳ, προέφθασαν πρὸς ὄρθρον, ἐκζητοῦσαι ὡς ἡμέραν,
Μυροφόροι κόραι, καὶ πρὸς ἀλλήλας ἐβόων· Ὡ φίλαι, δεῦτε τοῖς ἀρώμασιν ὑπαλείψωμεν,
Σῶμα ζωηφόρον καὶ τεθαμμένον, σάρκα ἀνιστῶσαν τὸν παραπεσόντα Ἀδὰμ κείμενον ἐν τῷ μνήματι, ἄγωμεν,
σπεύσωμεν, ὥσπερ οἱ Μάγοι, καὶ προσκυνήσωμεν, καὶ προσκομίσωμεν τὰ μύρα ὡς δῶρα τῷ μὴ ἐν σπαργάνοις,
ἀλλ' ἐν σινδόνι ἐνειλημένῳ, καὶ κλαύσωμεν, καὶ κράξωμεν· Ὡ Δέσποτα ἐξεγέρθητι, ὁ τοῖς πεσοῦσι παρέχων ἀνάστασιν.]

[Συναξάριον]

[Τῇ ἀγίᾳ καὶ μεγάλη Κυριακῇ τοῦ Πάσχα, αὐτὴν τὴν ζωηφόρον Ἀνάστασιν ἐορτάζομεν τοῦ Κυρίου, καὶ Θεοῦ καὶ Σωτῆρος ἡμῶν Ἰησοῦ Χριστοῦ.]

[Στίχοι]

[Χριστὸς κατελθὼν πρὸς πύλην Ἀδου μόνος]

[Λαβὼν ἀνῆλθε πολλὰ τῆς νίκης σκῦλα.]

[Ἰκός:]

[Еже прѣжде со́лнца Со́лнце заше́дшее иногд́а во гроб, предвари́ша ко́ у́тру, и́щущия́ я́ко дне миронóсицы дѣвы,
и друѓа ко друзѣ́й вопи́яху: о друѓини! Приид́ите, вон́ями помáжем тѣло живонóсное и погребѣ́ное,
плоть воскресѣ́вшаго пáдшаго Ада́ма, лежа́щую во грóбе. И́дем, потщ́имся я́коже волсвѣ́й, и поклонѣ́мся, и принесѣ́м мѣ́ра я́ко дáры,
не в пеленáх, но в плащани́це обвѣ́тому, и плáчим, и возопи́йм: о Влады́ко, востáни, пáдшим пода́й воскресѣ́ние.]

[Синаксáрь, во свyтýю и вел́икyю Недѣ́лю Пáсхи.]

[Стих́и: Христóс сше́д к борьбѣ́ áдове Ед́ин.]

[Мно́гия взѣ́м побѣ́ды коры́сти, взы́де.]

[Во свyтýю и вел́икyю Недѣ́лю Пáсхи, сáмoе живонóсное Воскресѣ́ние прáзднуем Гóспода Бóга и Спáса нáшего Иисýса Христá]

CSL	GR	1	2	3	4	5	6	7
1	[0.393678	0.04205396	-0.24479382	-0.05323773	-0.24479382	0.0911667	-0.01963074]
2	[-0.01602886	0.34094387	-0.13031456	0.22283232	-0.13031456	0.1565712	0.11568692]	
3	[0.07481634	0.21875925	-0.17499971	0.5257409	-0.17499971	0.23422936	-0.05796702]
4	[0.2000877	0.11701456	-0.23621589	0.16487566	-0.23621589	0.26789212	0.02254544]
5	[0.00616331	0.0261609	0.08375084	-0.0773346	0.08375084	0.03489467	0.359214]
6	[0.03978096	0.32326674	-0.27889502	0.7659529	-0.27889502	0.27096146	0.00304512]]

3. [Συναξάριον]

4. [Τῇ ἀγίᾳ καὶ μεγάλῃ Κυριακῇ τοῦ Πάσχα, αὐτὴν τὴν ζωηφόρον Ἀνάστασιν ἐορτάζομεν τοῦ Κυρίου, καὶ Θεοῦ καὶ Σωτῆρος ἡμῶν Ἰησοῦ Χριστοῦ.]

5. [Στίχοι]

6. [Χριστὸς κατελθὼν πρὸς πύλην Ἰδοῦ μόνος]

7. [Λαβὼν ἀνῆλθε πολλὰ τῆς νίκης σκῦλα.]

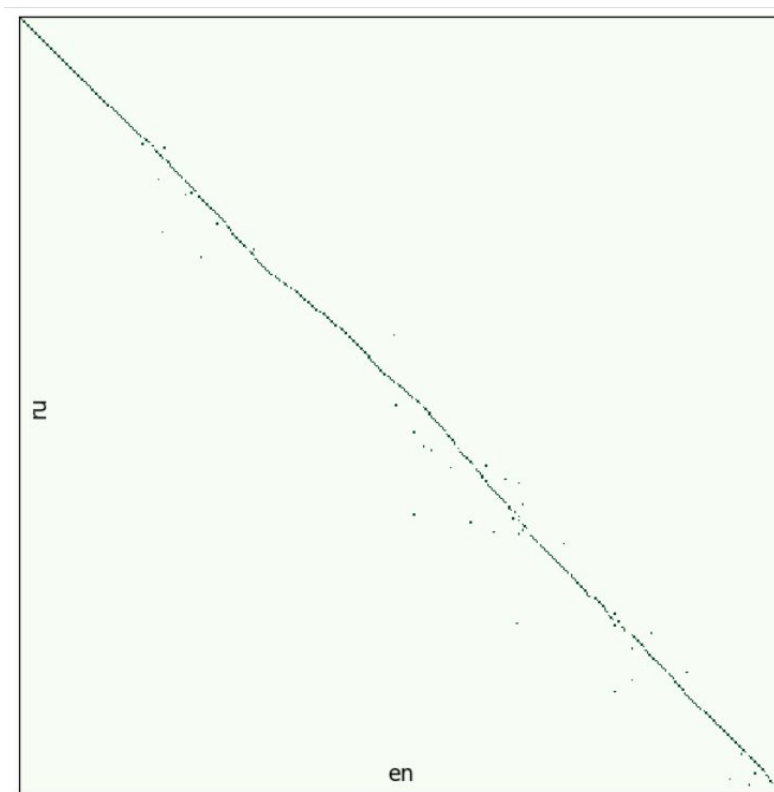
3. [Синаксáрь, во святую и великую Неделю Пáсхи.]

4. [Στιχί: Χριστὸς сшéd к борьбé áдове Ед́ин.]

5. [Мнóгия взém побéды корýсти, взýде.]

6. [Во святую и великую Неделю Пáсхи, сáмое живонóсное Воскресéние прázднуем Гóспода Бóга и Спáса нáшего Иисýса Христá]

Дистрибутивная семантика и выравнивание текстов



Спасибо за внимание!