

# Корпусная лингвистика

*Corpus linguistics*

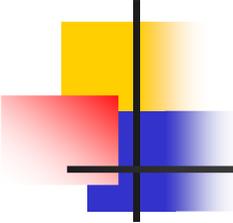


Петрозаводский государственный  
университет



Крижановский Андрей Анатольевич

andrew.krizhanovsky  gmail.com 1

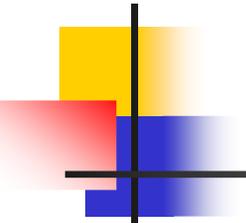


# Основные понятия

---

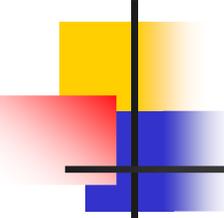
**Корпусная лингвистика** – это наука о создании и использовании текстовых (лингвистических) **корпусов**, возникшая вследствие растущих потребностей лингвистики во внедрении компьютерных технологий для работы с большими массивами языковых данных.

# What is Corpus Linguistics?



---

**Corpus Linguistics** is the study of language/linguistic phenomena through the analysis of data obtained from a corpus.



# Лингвистический корпус (1)

---

- совокупность **ТЕКСТОВ**,
- собранных в соответствии с определёнными принципами,
- размеченных по определённому стандарту
- и обеспеченных специализированной **ПОИСКОВОЙ СИСТЕМОЙ**.



## Лингвистический корпус (2)

---

Под текстами в этом случае понимаются не только продукты **письменного** языка (газетные статьи, романы, письма, электронные сообщения, дневники и т.п.), но и **устные** высказывания (доклады, радио- и телепередачи, телефонные разговоры и т.п.).

# Why to use a *corpus*?

- Intuition alone is not enough
  - Is “*starting*” always replaceable by “*beginning*”?
  - Is it only “*time*” that is “*immemorial*”?
  - “*think of*” vs. “*think about*”
- Native speaker intuition is unreliable
  - provides no information on frequency of occurrence
  - “***head***” => body part - Is this the most used sense?
- Один раз создать корпус и многократно применять его для решения различных лингвистических задач.

# Text vs. Corpus

(Tognini-Bonelli 2001: 3)

TEXT	CORPUS
Read whole	Read fragmented
Read horizontally	Read vertically
Read for content	Read for formal patterning
Read as a unique event	Read for repeated events
Read as an individual act of will	Read as a sample of social practice
Coherent communicative event	Not a coherent communicative event

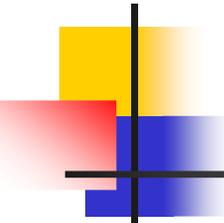
# Text vs. Corpus

From time to time there is also the need for high quality information to support particular initiatives, such as the (successful) application for accreditation. Some progress has been made in recording data on the Polytechnic 's rooms and buildings, and on the teaching space requirements of individual courses. These data are analysed, along with the database on course details and students ' course and module registrations, using the methodology in DES Design Note 44. Ad hoc reports are an essential part of any system that aspires not merely to process data routinely but to permit management information to be creamed off the top.

## N

## Concordance

13 ement system. They can choose whether to enter **data** themselves or to use the data preparation se  
14 individuals is recorded. As well as student-related **data**, details on courses, on modules and their  
15 student on the module is entered, for subsequent **data** processing by the registry. In addition  
16 s detailed record-keeping and effortless access to **data**. All of this the system provides.  
17 nd passed to the Registry who, together with the **data** preparation service, enter and verify approxi  
18 considerable use of student management system **data** processing. The marksheet for each module  
19 ons that can assist their work. Individual student **data** are available to assist counselling. Registry  
20 tion. Some progress has been made in recording **data** on the Polytechnic 's rooms and buildings, a  
21 onitors to allow the whole Committee to view such **data**. A detailed analysis of the performance of  
22 space requirements of individual courses. These **data** are analysed, along with the database on c  
23 easy to record in a computer (unless complicated **data** structures are used) and are even harder to



# Тексты, входящие в состав корпусов

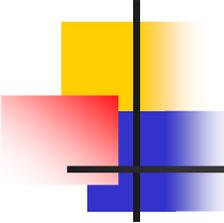
---

- отобраны исходя из определенных принципов,
- специально подготовлены и размечены,
  - машиночитаемый формат + разметка
- с помощью специальных программ в них можно искать необходимые фрагменты текста по заданным параметрам.

# Классификация корпусов: критерии

Характеристики анализа и сравнения корпусов:

- Тип/формат данных
- Язык текстов
- «Параллельность»
- «Литературность»
- Специфичность
- Жанр
- Доступность
- Назначение
- Динамичность
- Разметка
- Характер разметки
- Объем текстов
- Хронологический аспект
- «Общность»
- Структура



# Основные понятия

---

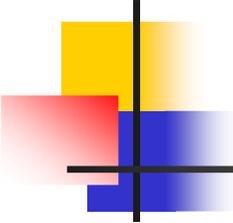
## Типы корпусов:

1. Противопоставление корпусов, относящихся ко **всему языку**, корпусам, относящимся к какому-либо **подъязыку** (жанр, стиль, язык определенной возрастной или социальной группы, язык писателя или ученого и т.п.);

2. Разделение корпусов по типу лингвистической разметки.

Несмотря на наличие множества типов разметки, большинство реально существующих корпусов относится к корпусам **морфологического** либо **синтаксического** типа (treebanks, «банки синтаксических структур»).

<b>Критерий</b>	<b>Типы корпусов</b>
Формат текста	Электронный / Неэлектронный
Полнота текстов	Полнотекстовый / Выборочный
Завершенность корпуса	Статичный / Пополняемый
Средство реализации языка	Корпус письменной речи Корпус устной речи
Языковые разновидности	Корпус стандартного языка Корпус молодежного языка Корпус экономического / компьютерного языка
Временной параметр	Корпус современного языка Корпус исторического языка
Количество языков	Одноязычный / Многоязычный

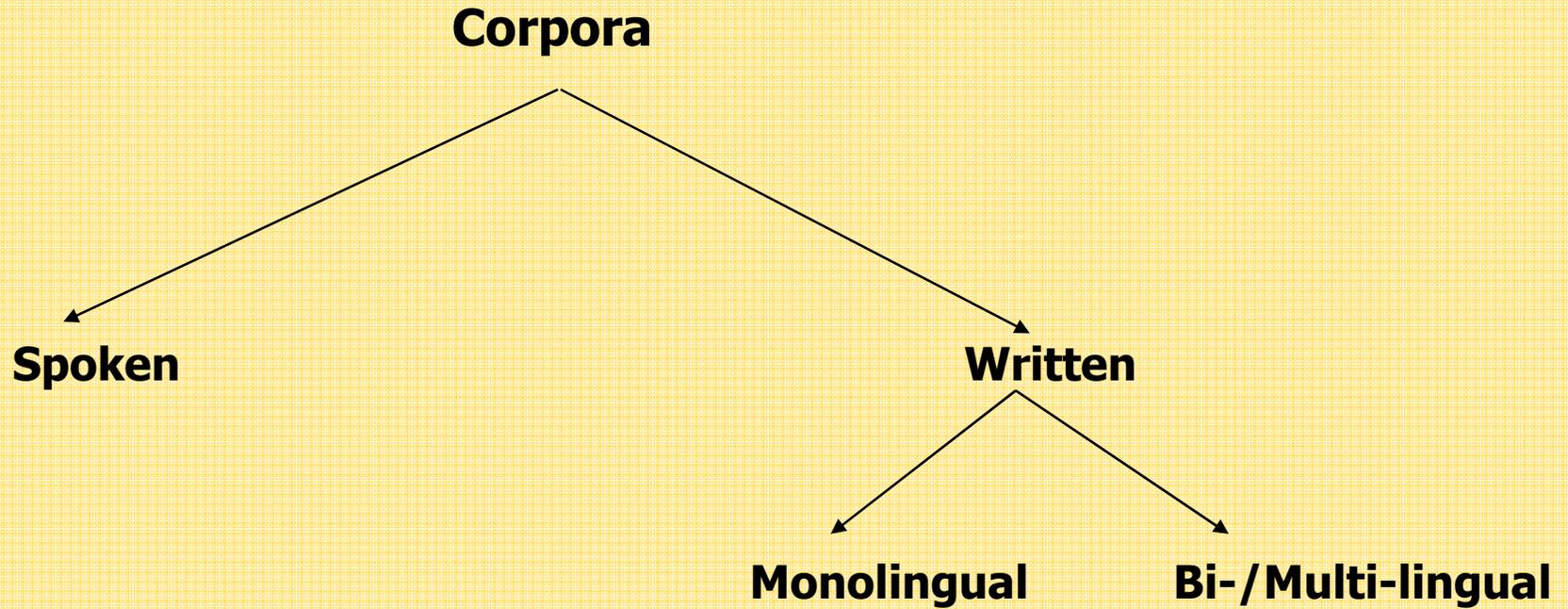


# Types of corpora

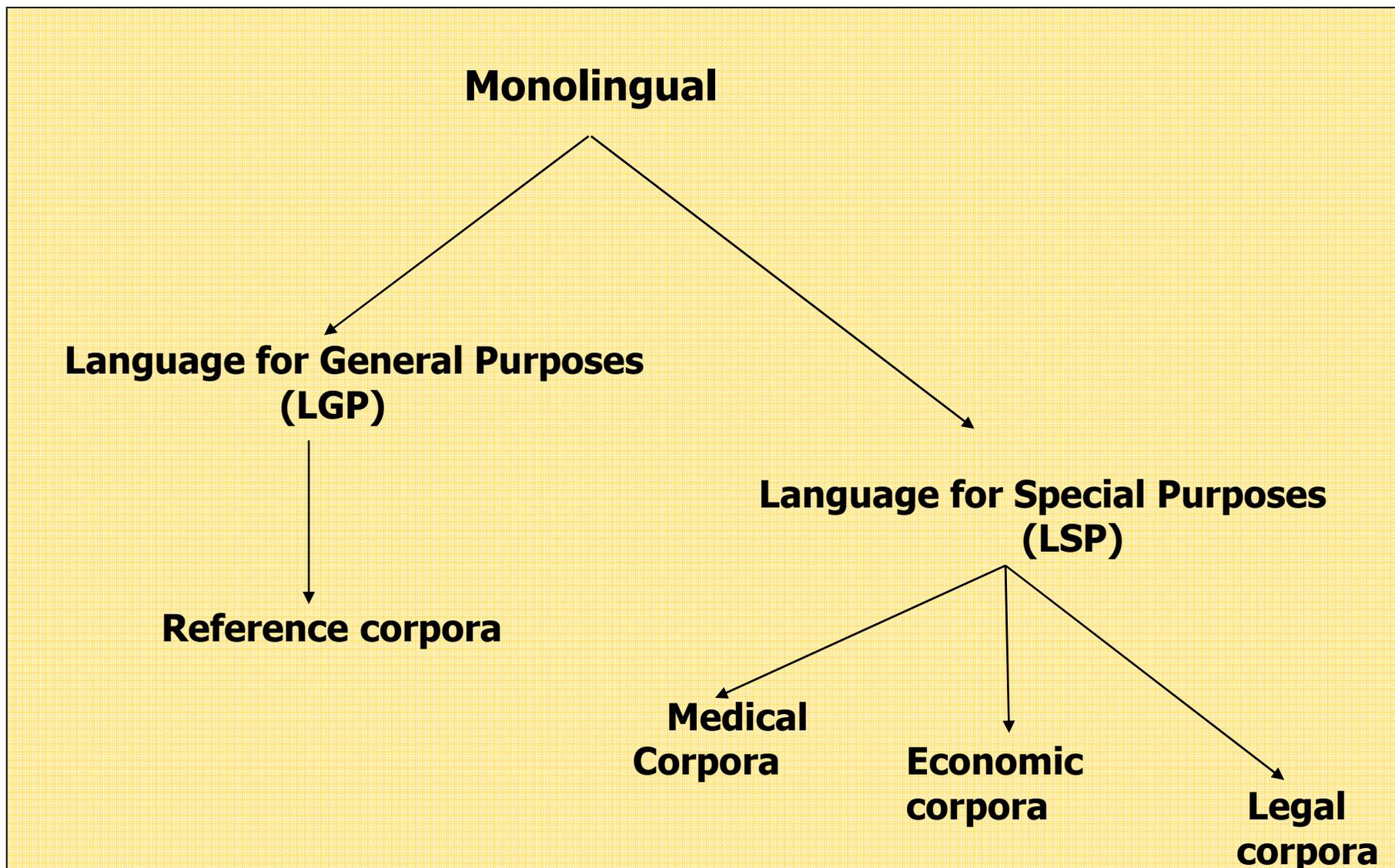
---

- spoken vs. written
- monolingual vs. bi/multilingual
- parallel vs. comparable corpora  
(translation corpora)
- general language purpose vs.  
specialised  
language purpose
- diachronic vs. synchronic
- plain text vs. annotated (tagged) text

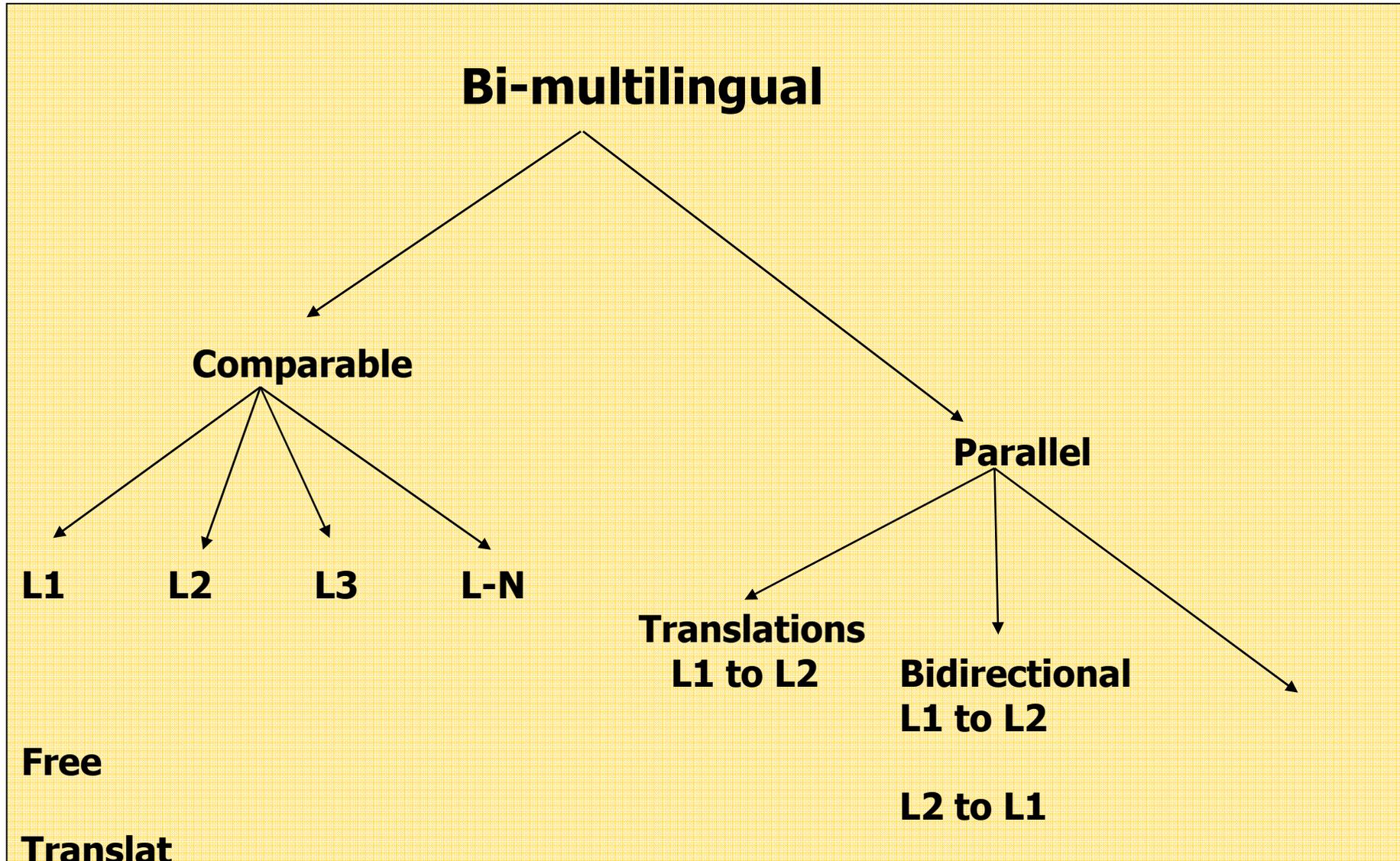
# Types of corpora



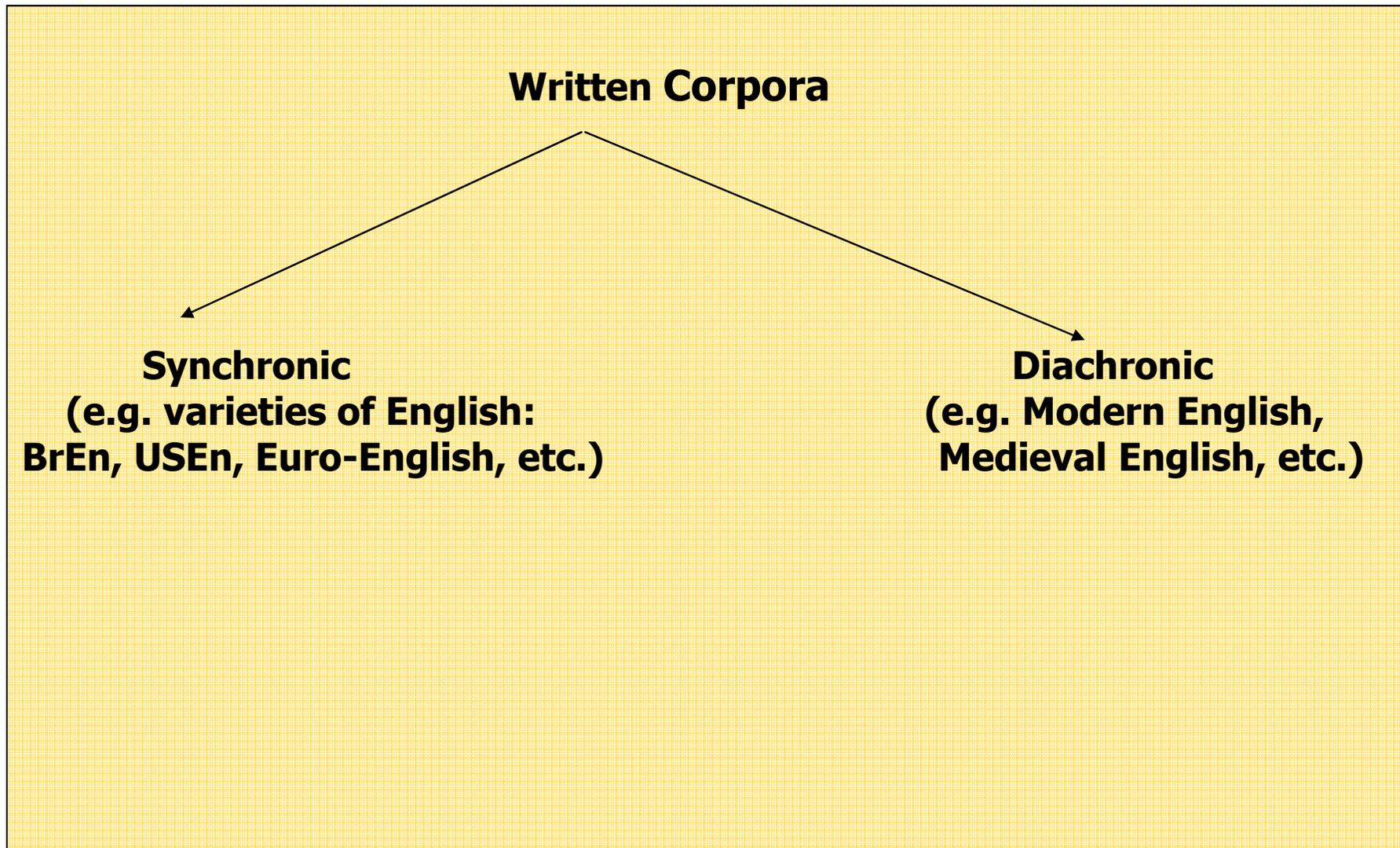
# Types of corpora

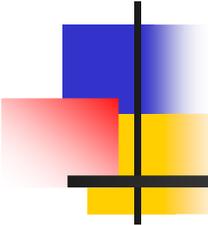


# Types of corpora



# Types of corpora





# Корпусы(а) в Интернете

---

<p>Национальный корпус русского языка  <a href="http://ruscorpora.ru">http://ruscorpora.ru</a></p>	<p>&gt;300 млн  слов</p>
<p>Открытый корпус русского языка  <a href="http://opencorpora.org">http://opencorpora.org</a></p>	<p>&gt;921 тыс.  слов</p>
<p>Компьютерный корпус текстов русских газет конца  XX-го века  <a href="http://www.philol.msu.ru/~lex/corpus">http://www.philol.msu.ru/~lex/corpus</a></p>	<p>200 тыс. слов</p>
<p>Корпус русского языка ХАНКО (Хельсинский  университет)  <a href="http://www.ling.helsinki.fi/projects/hanco/">http://www.ling.helsinki.fi/projects/hanco/</a></p>	<p>100 тыс. слов  Ручная  морфологическая  разметка</p>
<p>Корпуса русских текстов на сайте Университета в  Лидсе, Великобритания  <a href="http://corpus.leeds.ac.uk">http://corpus.leeds.ac.uk</a></p>	
<p>Русские корпуса Тюбингенского Университета  <a href="http://www.sfb441.uni-tuebingen.de/b1/en/korpora.html">http://www.sfb441.uni-tuebingen.de/b1/en/korpora.html</a></p>	
<p>Словарь-корпус языка А.С. Грибоедова  <a href="http://feb-web.ru/feb/concord/abc/">http://feb-web.ru/feb/concord/abc/</a></p>	<p>120 тыс. слов</p>

<p>Уппсальский корпус русских текстов Доступен для поиска на сайте <b><a href="http://www.sfb441.uni-tuebingen.de/b1/en/korpora.html">http://www.sfb441.uni-tuebingen.de/b1/en/korpora.html</a></b></p>	<p>1 млн слов 600 текстов (публицистика 1985-1989; литературные произведения 1960-1988)</p>
<p>Банк английского языка (Bank of English) <b><a href="http://www.collins.co.uk/books.aspx?group=153">http://www.collins.co.uk/books.aspx?group=153</a></b> Свободный доступ: <b><a href="http://www.collins.co.uk/Corpus/CorpusSearch.aspx">http://www.collins.co.uk/Corpus/CorpusSearch.aspx</a></b></p>	<p>524 млн слов, 56 млн в свободном доступе (The Collins Wordbanks <i>Online</i> English corpus: 36 млн – брит. англ., 10 млн – амер. англ., 10 млн – брит. разговорн. англ.)</p>
<p>Британский национальный корпус <b><a href="http://www.natcorp.ox.ac.uk/">http://www.natcorp.ox.ac.uk/</a></b> или <b><a href="http://sara.natcorp.ox.ac.uk/">http://sara.natcorp.ox.ac.uk/</a></b></p>	<p>100 млн слов Корпусные менеджеры SARA и XAIRA (<b><a href="http://www.xaira.org">http://www.xaira.org</a></b>)</p>
<p>Венгерский национальный корпус <b><a href="http://corpus.nytud.hu/mnsz/">http://corpus.nytud.hu/mnsz/</a></b></p>	<p>100 млн слов</p>

<p>Корпус испанского языка (исторический)  <a href="http://www.corpusdelespanol.org/">http://www.corpusdelespanol.org/</a></p>	<p>100 млн слов, тексты 13–20 вв.  Создан в Иллинойском университете, США</p>
<p>Корпус современного датского языка  <a href="http://www.korpus2000.dk/">http://www.korpus2000.dk/</a></p>	<p>50 млн слов  Тексты 1998–2002 гг.</p>
<p>Корпус современного итальянского языка  CORIS/CODIS  <a href="http://www.cilta.unibo.it/ricerca.htm">http://www.cilta.unibo.it/ricerca.htm</a></p>	<p>100 млн слов</p>
<p>Корпус современного китайского языка  (LIVAC Synchronous Corpus)  <a href="http://www.rcl.cityu.edu.hk/livac/">http://www.rcl.cityu.edu.hk/livac/</a></p>	<p>720 млн слов  (150 млн иероглифов)</p>
<p>Мангеймский корпус немецкого языка  (Institut für Deutsche Sprache, Mannheim,  Germany)  <a href="http://corpora.ids-mannheim.de/ccdb/">http://corpora.ids-mannheim.de/ccdb/</a></p>	<p>1610 млн слов  Корпусный менеджер COSMAS</p>

Польский национальный корпус <a href="http://korpus.ia.uni.lodz.pl/">http://korpus.ia.uni.lodz.pl/</a>	93 млн слов
Словацкий национальный корпус <a href="http://korpus.juls.savba.sk">http://korpus.juls.savba.sk</a>	180 млн слов Используется корпусный менеджер Manatee/Bonito
Хорватский национальный корпус <a href="http://www.hnk.ffzg.hr/">http://www.hnk.ffzg.hr/</a>	53 млн слов Корпусный менеджер Manatee/Bonito
Чешский национальный корпус <a href="http://ucnk.ff.cuni.cz">http://ucnk.ff.cuni.cz</a>	100 млн слов + 100 млн нового корпуса современной лексики Корпусный менеджер Manatee/Bonito
Эстонский корпус <a href="http://www.cl.ut.ee/korpused/baaskorpus/">http://www.cl.ut.ee/korpused/baaskorpus/</a>	

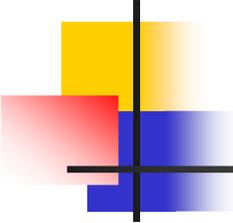


## **Национальный корпус русского языка**

представляет данный язык на определенном этапе его существования и во всём многообразии жанров, стилей, территориальных и социальных вариантов.

Образовательный портал Национального корпуса русского языка:

**<http://studiorum.ruscorpora.ru/>**

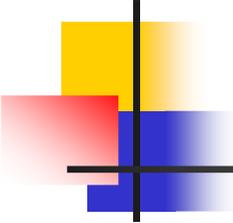


## НКРЯ (2)

---

Национальный корпус имеет две важные **особенности**:

1. Он характеризуется **представительностью**, или сбалансированным составом текстов.
2. Корпус содержит особую дополнительную информацию о свойствах входящих в него текстов (так называемую **разметку**, или **аннотацию**).

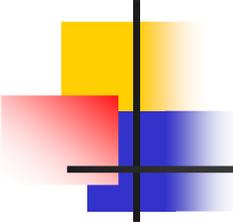


# Основные понятия

---

**Конкорданс** – это список всех употреблений заданного языкового выражения (например, слова) **в контексте**, возможно, со ссылками на источник.

Существуют специальные программы составления конкордансов по некоторому корпусу текстов, так называемые **конкордансеры**.



# Основные понятия

---

Итак, корпуса используются, прежде всего, при

- исследовании различных языковых разновидностей; проверка лингвистических теорий;
- составлении словарей, грамматических справочников и т.п.;
- изучении и преподавании иностранных языков (составление учебных пособий, упражнений, доступ к аутентичным материалам);
- переводе (проверка сочетаемости слов, стилистических особенностей);
- социолингвистике: сравнение языков разных социальных групп;

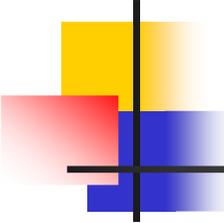
# Lexicography / Terminology (wikipedia.org)

General **lexicography** focuses on the design, compilation, use and evaluation of **general dictionaries**, i.e. dictionaries that provide a description of the language in general use.

**Terminology**, in its general sense, simply refers to the usage and study of terms, that is to say words and compound words generally used in **specific contexts**.

# Lexicography and corpora

- Corpus-based lexicography started in England
- Corpus provides authentic uses of language
- **Extract samples (concordance) to identify different senses**
- Word Frequency information
- **Help identify** collocation (1), set phrase (2)
  - 1) Фразеологическое сочетание, коллокация
    - *ставить условия, вносить предложения*
  - 2) Фразеологическое выражение
    - *пословицы, афоризмы, речевые штампы*
      - всего хорошего, до новых встреч
- Most English dictionaries are now corpus-based. Oxford, Collins, Longman, Cambridge, Macmillan,



# Источники литературы

1. Захаров В.П. Обзор корпусов. Презентация. – Режим доступа: [download.yandex.ru/class/zakharov/CL\\_L9.ppt](http://download.yandex.ru/class/zakharov/CL_L9.ppt)
  2. Образовательный портал Национального корпуса русского языка. – Режим доступа: [http://studiorum.ruscorpora.ru/index.php?option=com\\_content&view=article&id=241&Itemid=48](http://studiorum.ruscorpora.ru/index.php?option=com_content&view=article&id=241&Itemid=48)
  3. Подлеская В.И. Современные компьютерные методы в изучении и преподавании лингвистических дисциплин: корпусная лингвистика. – Режим доступа: [http://zhangbyrzhan.ucoz.ru/publ/metodika/inostrannyj\\_jazyk/sovremennye\\_kompjuternye\\_metody\\_v\\_izuchenii\\_i\\_prepodavanii\\_lingvisticheskikh\\_disciplin/12-1-0-27](http://zhangbyrzhan.ucoz.ru/publ/metodika/inostrannyj_jazyk/sovremennye_kompjuternye_metody_v_izuchenii_i_prepodavanii_lingvisticheskikh_disciplin/12-1-0-27)
  4. Портал «Национальный корпус русского языка». – Режим доступа: <http://www.ruscorpora.ru/>
  5. Портал «Фонд знаний ЛОМОНОСОВ». Энциклопедия. Статья «Конкорданс». – Режим доступа: <http://www.lomonosov-fund.ru/enc/ru/encyclopedia:0127200>
  6. Scherer C. Korpuslinguistik. – Universitätsverlag WINTER Heidelberg. – 2006. – 98 S.
- Презентации:
1. Электронные корпуса. (безымянная презентация)  
[http://www.slideshare.net/anna\\_pal/ss-13040000](http://www.slideshare.net/anna_pal/ss-13040000)
  2. Corpus linguistics: a general introduction (who is author?)  
[http://www.lingue.uniba.it/dag/pagine/personale/falco/Corpus%20linguistics\\_introduction\\_Students.ppt](http://www.lingue.uniba.it/dag/pagine/personale/falco/Corpus%20linguistics_introduction_Students.ppt)