

Гибкий Павел

Корпусы текстов КИТАЙСКОГО ЯЗЫКА



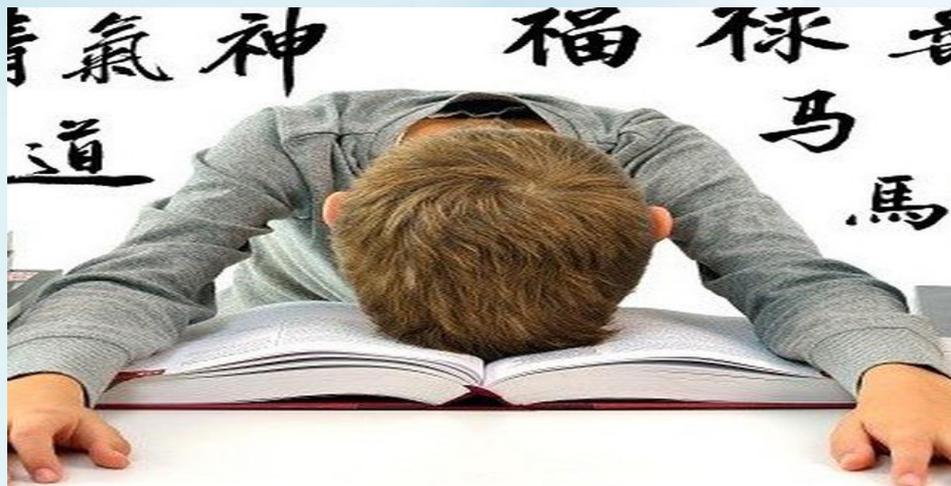
«Конец 1980 – середина 1990-х гг., создаются корпусы текстов на национальных языках в разных странах, в том числе и в Китае» [1].



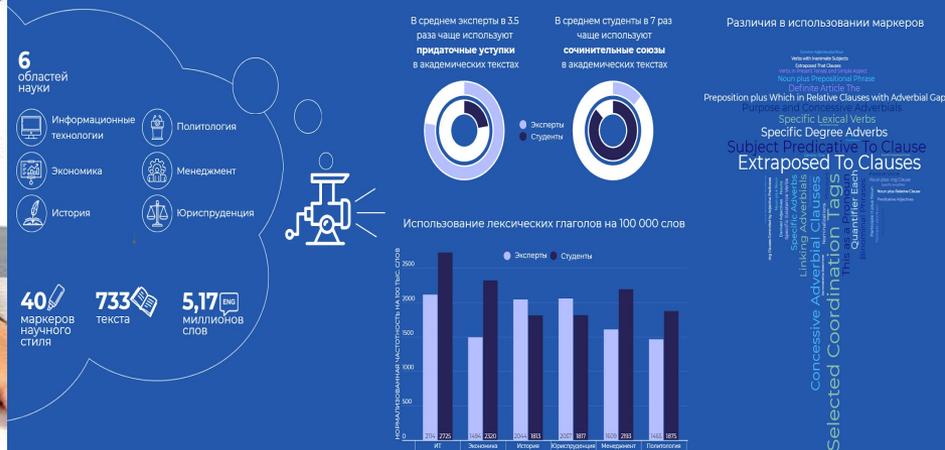
人民日报

Предвестник корпусов китайского языка

Собрание китайских текстов для исследования частотности 语体文应用字汇 («Сборник текстов для изучения единиц разговорного стиля языка») [3].



Академические тексты: эксперты vs. студенты



Первый китайский лингвистический
корпус

«人民日报»标注语料库 (
корпус газеты
«Жэньминьжибао» (1999
г. [3]).



Современные корпуса китайского языка

- «Лингвистический корпус китайского языка Пекинского университета языка и культуры (BCC).
- Center for Chinese Linguistics (CCL) .
- Chinese Corpus online (语料库在线, языковые материалы с 1919 года)» [2].
- НКРЯ (Русско-китайский параллельный корпус Национального корпуса русского



ВСС

Крупнейший корпус китайского
языка **в мире** (15 млрд
иероглифов).

A globe with a grid overlay, showing the continents. The text 'ВСС:' is written in large, red, pixelated characters across the globe. The globe is slightly blurred, and the grid lines are prominent.

ВСС:

Chinese Corpus online

Программы автоматической
сегментации текстов,
частеречной разметки слов,
подсчета частотности слов и
разметки пиньиня (100 млн

4



ССЛ

Корпус современного,
древнекитайского
языков, китайско-английский
корпус (500 млн иероглифов).



НКРЯ

- самый большой
открытый параллельный
корпус русского и
китайского языков.



НАЦИОНАЛЬНЫЙ КОРПУС
РУССКОГО
ЯЗЫКА

Что такое НКРЯ?

КОРПУС РУССКОГО ЯЗЫКА — ЭТО
ИНФОРМАЦИОННО-СПРАВОЧНАЯ СИСТЕМА,
ОСНОВАННАЯ НА СОБРАНИИ РУССКИХ
ТЕКСТОВ В ЭЛЕКТРОННОЙ ФОРМЕ. ОБЩИЙ
ОБЪЁМ: БОЛЕЕ 600 МЛН СЛОВ.

Особенности китайских корпусов текстов

«Не все иероглифы китайского языка характеризуются **высокой встречаемостью** в текстах. Характерной особенностью иероглифов является их несоответствие **буквенно-словесным универсалиям**» [5].



мама папа

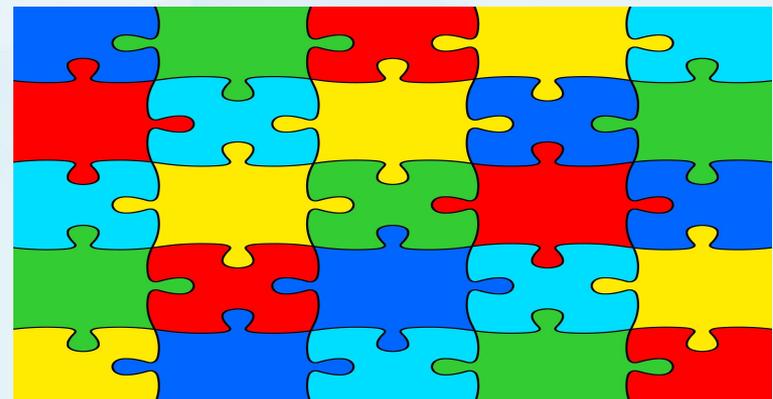
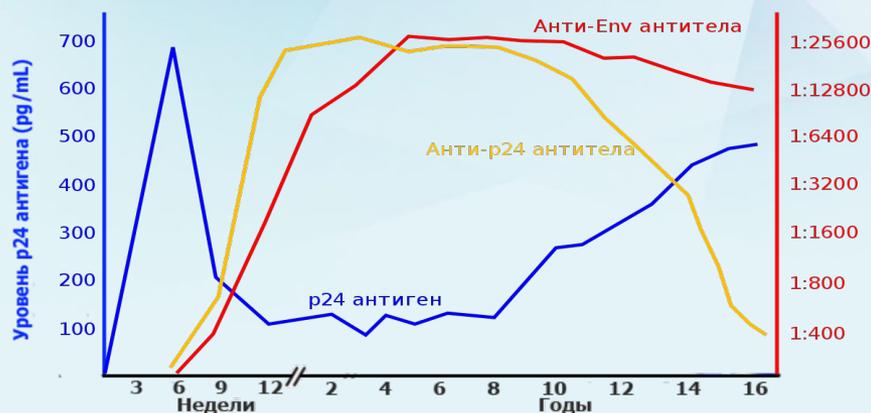
Проблемы ранних корпусов:

1. Большая часть данных вводилась вручную, небольшие размеры корпусов.



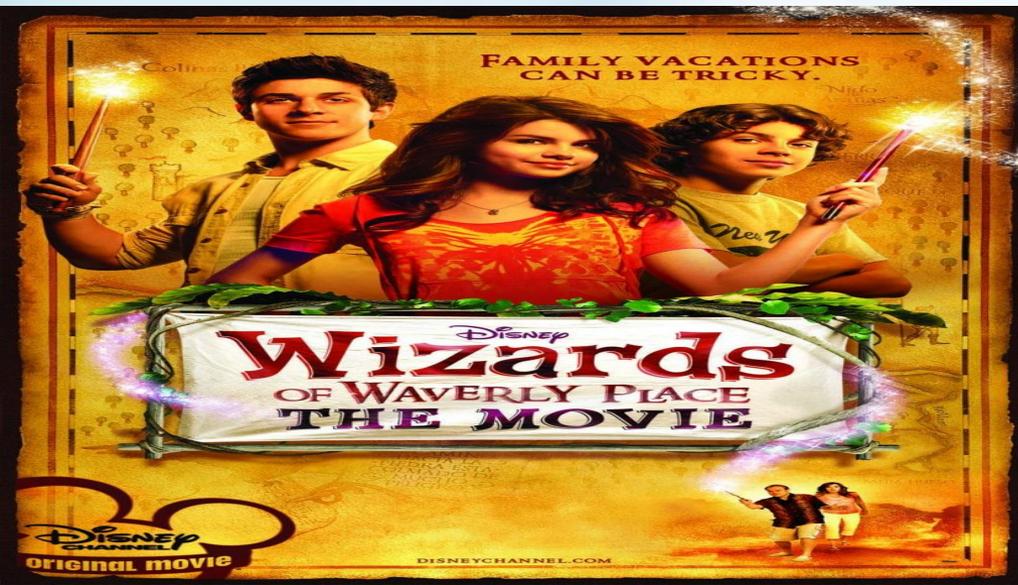
Проблемы ранних корпусов:

2. В силу использования различными корпусами разных методов автоматической сегментации получались разные результаты.



Решение проблем:

- 1) Разработка национального стандарта GB-13715 «Современная сегментация китайского слова ...» .
- 2) Составление первого масштабного корпуса китайского языка — 语料库在线 [4] (1991 год).
- 3) Составление CLL (783 463 175 знаков).



Список литературы

1. Баркович, А. А. Лингвистические корпусы китайского языка: функциональный аспект / А. А. Баркович, Ван Цин // Вестник МГЛУ. – 2015.. – . – Т. № 5, № (78). – С. 105 – 113.
2. Фэн, Юэ. Специфика корпусных исследований в современном китайском языкознании / Юэ Фэн, Ван Цин // Вестник МГЛУ. – 2020. – . – Т. 3, № 832. – С. 159 – 172.
3. 陈鹤琴 语体文应用字汇 = Сборник текстов для изучения языковых единиц разговорного стиля [Электронный ресурс]. – Режим доступа : http://book.ln.chaoxing.com/ebook/read_11378972.html. – Дата доступа : 25.04.2015.
4. 人民日报»标注语料库 = Размеченный корпус газеты «Жэньминьжибао» [Электронный ресурс]. – Режим доступа : <http://ling.cass.cn/yingyong/courses/corpusbase.htm>. – Дата доступа : 25.04.2015.
5. 字、词 – 现代汉语 = Соотношения между символами, словами и морфемами [Электронный ресурс]. – Режим доступа : <http://www.yyxx.sdu.edu.cn/chinese/wt/main04-03.htm>. – Дата доступа : 25.04.2015.