

НИС Культуромика

Борис Орехов

НИУ Высшая школа экономики

nevmenandr@gmail.com

19 октября 2018

- 1 Что такое «культуромика»
- 2 А что же с русским языком?
- 3 Что мы будем делать?

- 1 Что такое «культуромика»
- 2 А что же с русским языком?
- 3 Что мы будем делать?

Google решил создать искусственный интеллект и научить его на всех книгах, изданных человечеством.

Но книги для этого нужно оцифровать:

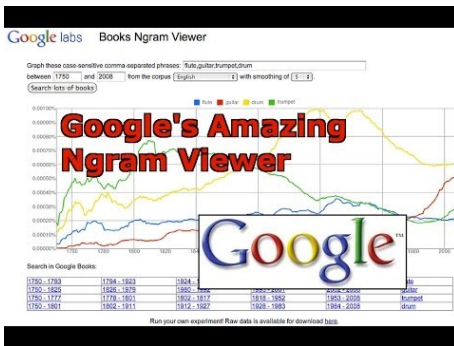


Google Ngram Viewer

В 2011 году Google был корпорацией добра и выложил собранные из книг текстовые данные в открытый доступ.

Из-за ограничений, связанных с авторскими правами тексты пришлось нарезать на n-граммы.

Кроме того, к ним был приделан поиск, в результате получился веб-сервис:

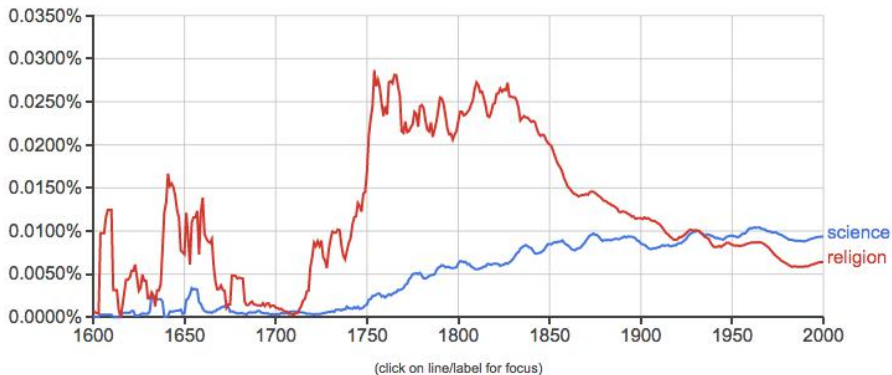


В 2011 году придумали «Культуромику».

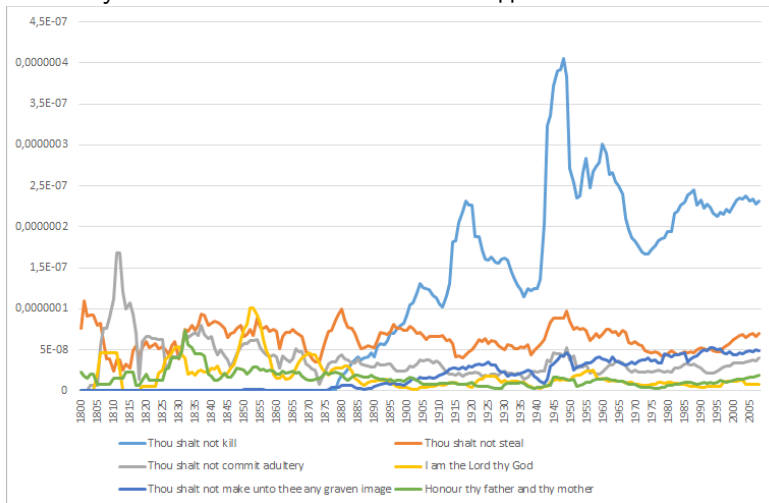
J.-B. Michel et al. Quantitative analysis of culture using millions of digitized books (2011)



Основная гипотеза: частотность употребления слов в книгах отражает важные культурные тренды



Частотность упоминания библейских заповедей



- 1 Что такое «культуромика»
- 2 А что же с русским языком?
- 3 Что мы будем делать?

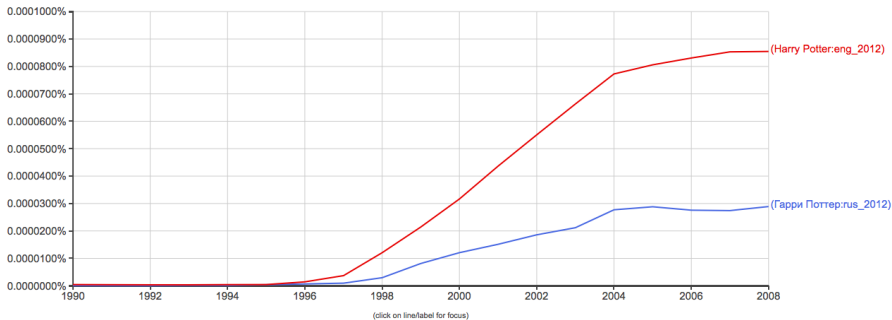
Книги на русском в данных есть

Гарри Поттер у нас и у них:

Google Books Ngram Viewer

Graph these comma-separated phrases: (Гарри Поттер:rus_2012),(Harry Potter:eng_2012) ☐ case-insensitive

between 1990 and 2008 from the corpus English with smoothing of 3 Search lots of books



Не распознана вся старая орфография

Ө

фрита

Ъ

ѣр

Ѣ

Ять

І

и десятеричное

Качество распознавания вообще страдает

- 1 Что такое «культуромика»
- 2 А что же с русским языком?
- 3 Что мы будем делать?

- Исправить орфографию в данных по дореволюционным книгам.
- Вывесить исправленные наборы ngrams в открытый доступ.
- Искать статистические выбросы и тренды автоматически:
 - Где пики и спады?
 - У каких слов и n-граммов отрицательная корреляция?
 - О чём не подумал исследователь, но что найдёт бездушная машина?
- Связывать слова и n-граммы с Wikidata и так автоматом искать связи.

Объяснение пика в связанных данных



Item [Discussion](#)

1980 Summer Olympics (Q8450)

Games of the XXII Olympiad, held in Moscow in 1980
Moscow 1980 | Games of the XXII Olympiad