

# Русский стилеметрически й датасет

уроки и возможности

Борис Орехов

## Назначение

Набор русскоязычных текстов,

подготовленный для

- стилеметрических экспериментов,
- машинного обучения,
- учебных занятий по компьютерному анализу текста



## Структура

Файлы сгруппированы по тематическим папкам



## Статус

Все тексты находятся в общественном достоянии (Public Domain) и доступны для исследований без ограничений.



## Временной охват

Произведения с 1763 по 1932 год (охватывает три века: XVIII, XIX и XX).



## Тщательный подбор

Все тексты снабжены метаданными и сопоставимы по объему внутри коллекции

## Состав датасета

### 01. Тексты

Более 320 документов

### 02. Слова

Совокупный объем — 16+  
миллионов слов

### 03. Авторы и подкорпуса

Тексты 94 авторов, распределенные  
по 25 подкорпусам

# Разделы датасета

Имена файлов соответствуют шаблону автор\_название.txt.  
Это позволяет пакету stylo (R) автоматически группировать  
тексты по классам при визуализации

**01**

## По авторам

Художественная  
проза XVIII–XX вв.,  
публицистика,  
научные труды

**02**

## Спецвыборк

Сравнение авторов  
по полу, возрасту,  
работе с  
псевдонимами  
(Чехов / Чехонте)

**03**

## Жанры и стили

Романтизм vs  
Реализм,  
Западники vs  
Славянофилы,  
стихи vs проза

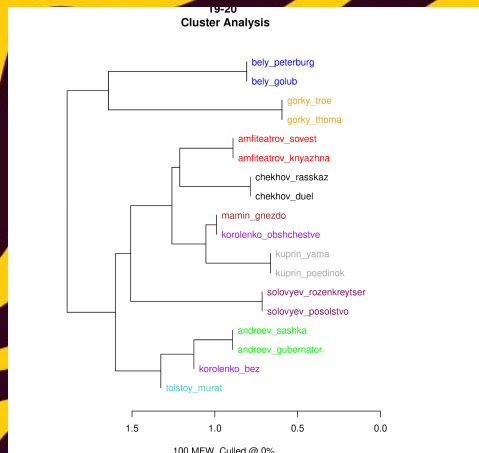
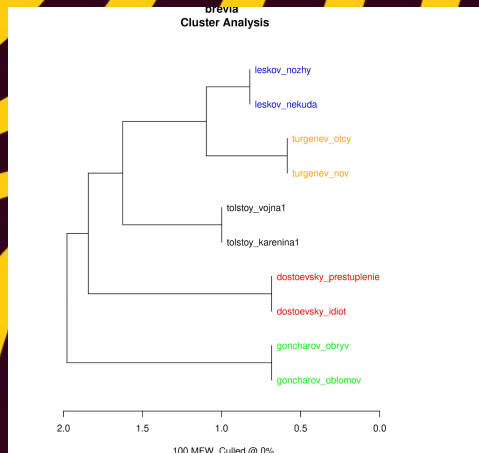
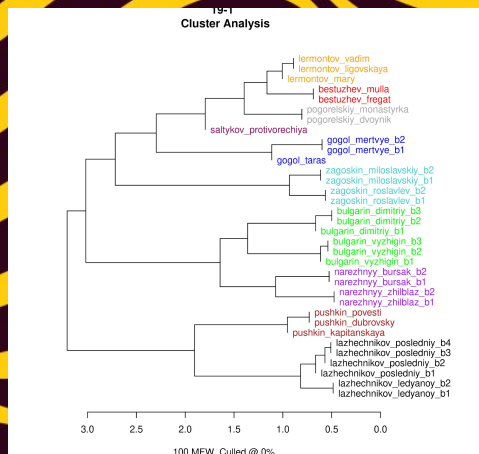
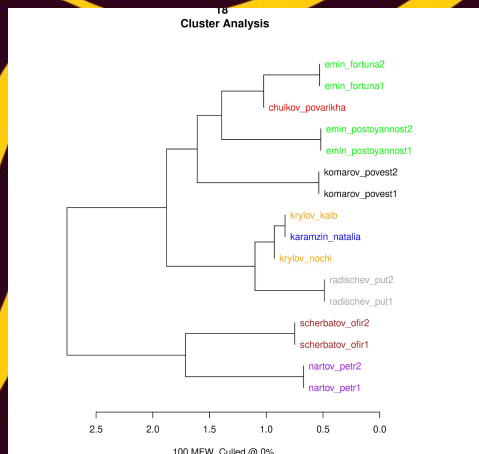
**04**

## Специфические

Речь персонажей и  
тексты,  
сгенерированные  
нейросетями



**Разведочный  
анализ: Delta,  
100 MFW**



# Авторы художественно й литературы

## Почти без сбоев

Не в свои кластеры  
попали только  
Карамзин и Короленко  
из 34 авторов

# Hacking stylometry with multiple voices: Imaginary writers can override authorial signal in Delta

Get access >

Daniil Skorinkin ✉, Boris Orekhov

*Digital Scholarship in the Humanities*, Volume 38, Issue 3, September 2023, 1266, <https://doi.org/10.1093/llc/fqad012>

Published: 08 April 2023

“ Cite   🔑 Permissions   ➦ Share ▼

## Abstract

## Автор и псевдоним

### Чехов и Чехонте

Тексты под псевдонимом и под реальной фамилией

### Псевдоним

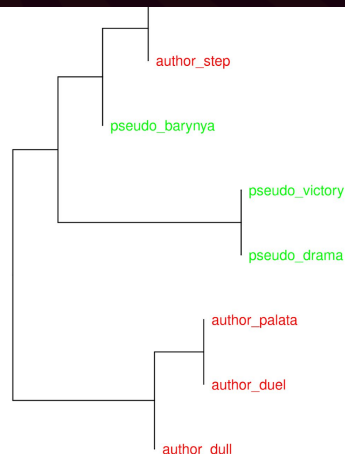
Драма на охоте, Цветы запоздалые, Ненужная победа

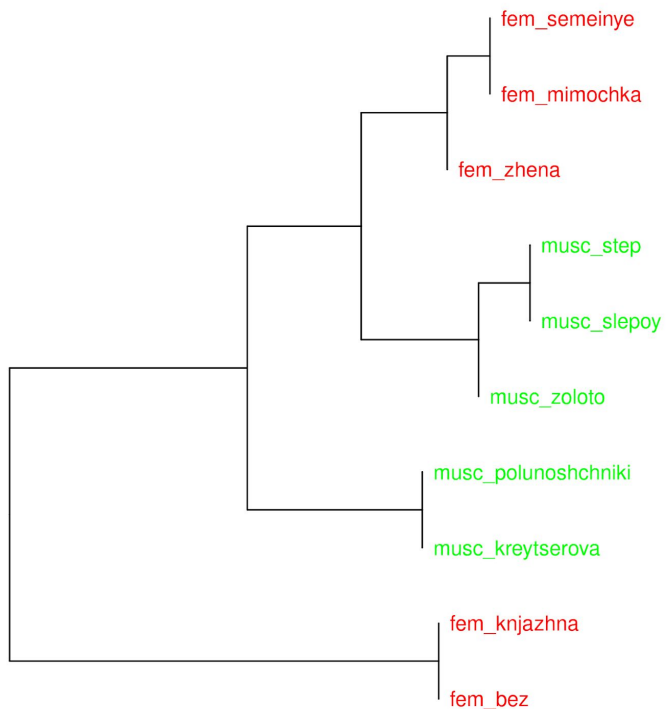
### Граница: начало 80-

Чехов переходит на практику публикации под своей фамилией

### Фамилия

Степь, Дуэль, Палата № 6





## Пол автора

Тексты 1880—  
1890-х годов

Не идеально, но  
обнадеживающе

# Поколения

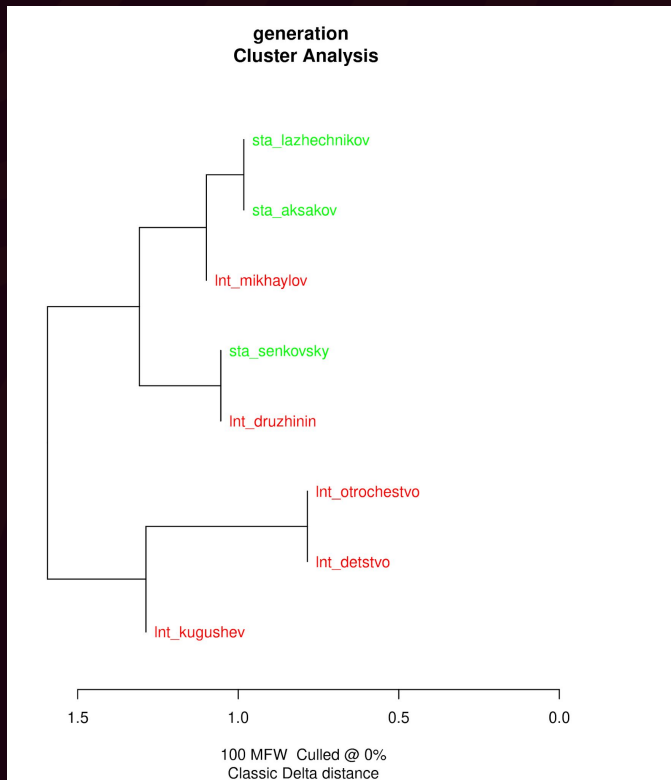


**Аксаков и Толстой**

**1856–1859**

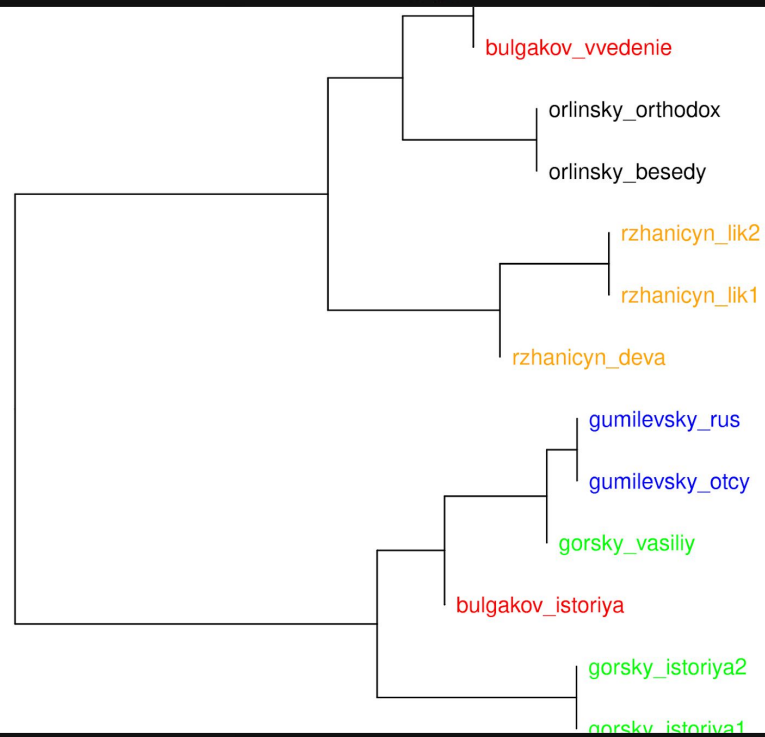
Лажечников, Сенковский  
(1790-е) vs. Дружинин и  
др. (1820-е)

«Семейная  
хроника», «Детство»

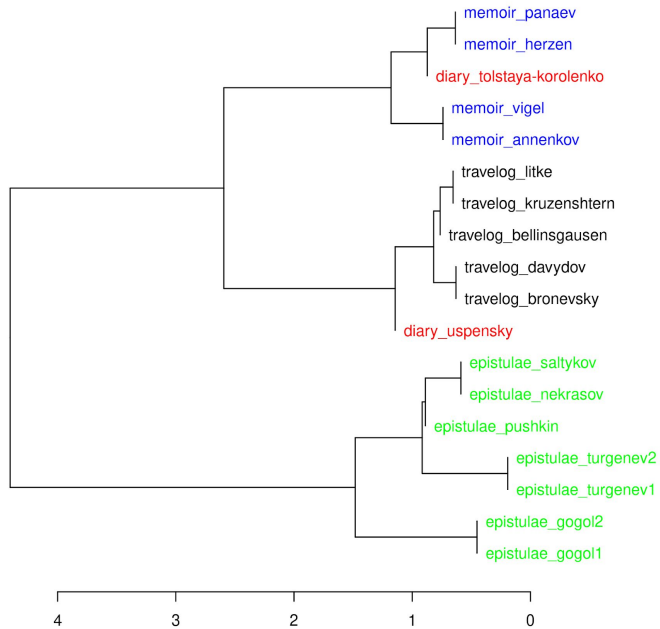


# Теология

Очень  
непоследовательные  
результаты  
Возможно, из-за  
библейских цитат



genre  
Cluster Analysis



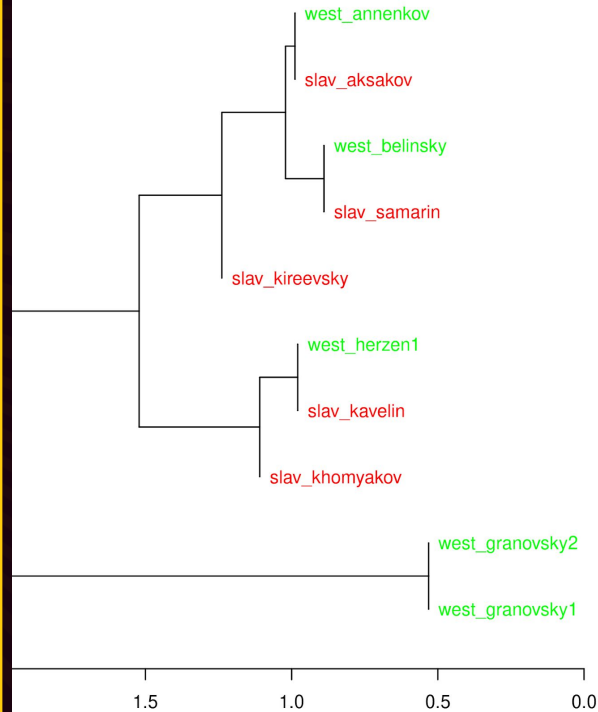
100 MFW Culled @ 0%  
Classic Delta distance

## Жанровый сигнал

Size matters

Дневники ~10 тыс слов; остальные от 80 до 200 тыс

### journalism Cluster Analysis



100 MFW Culled @ 0%  
Classic Delta distance

# Западники славянофилы



Идейные течения



Кажется, нет

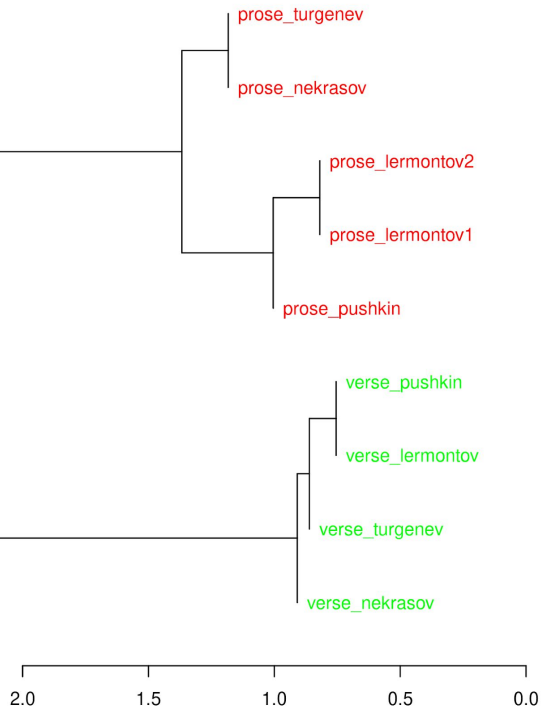
Умеет ли Delta отличать  
идеологию?

# Стих и проза



Самый явный сигнал

prose  
Cluster Analysis

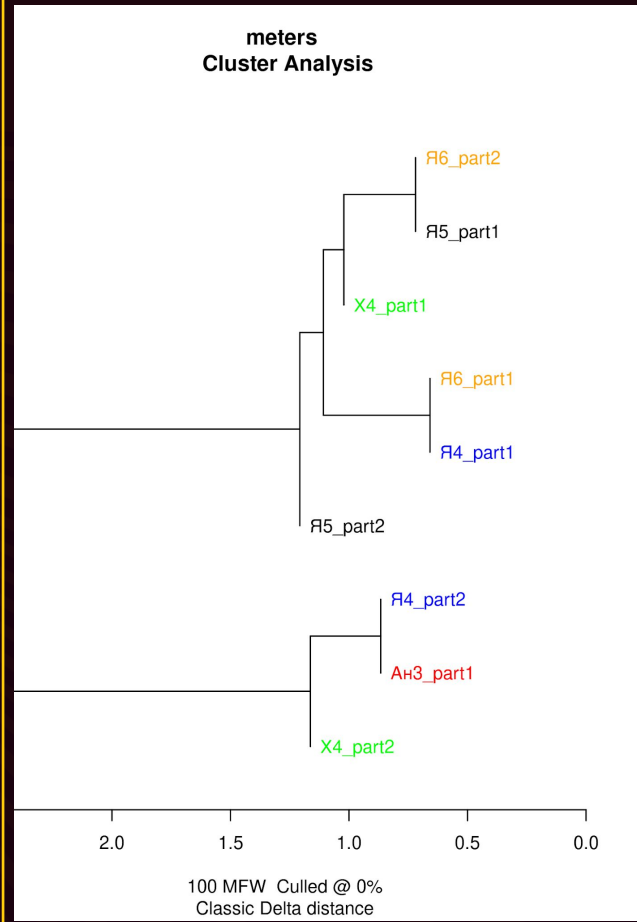


100 MFV Culled @ 0%  
Classic Delta distance

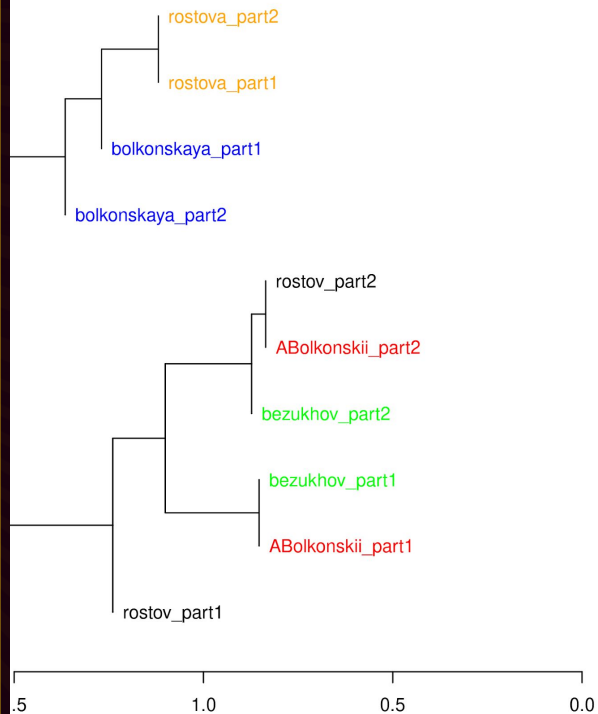
# Поэтические размеры



Не различаются



### tolstoy-voyna-i-mir Cluster Analysis



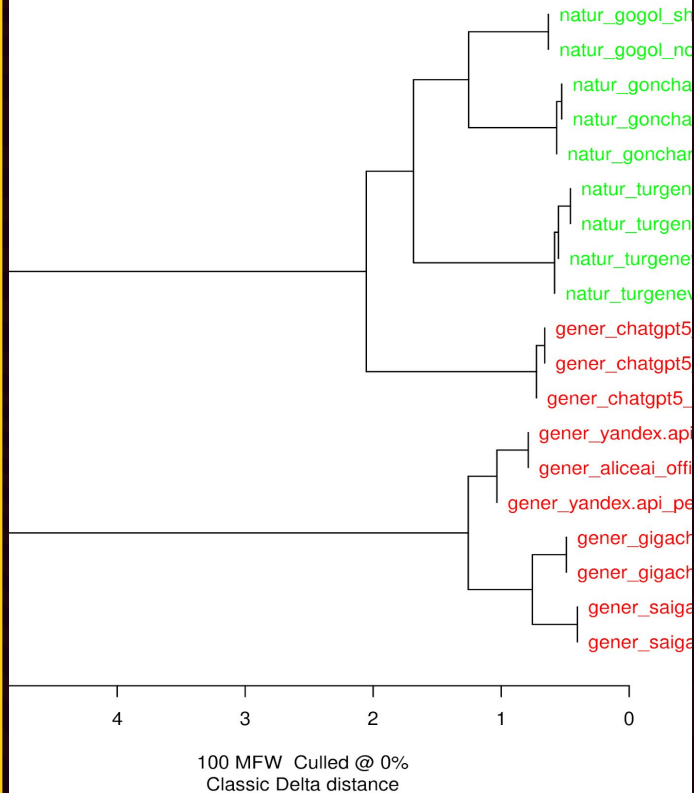
100 MFW Culled @ 0%  
Classic Delta distance

## Персонажи «Войны и мира»



К Д. А. Скоринкину есть вопросы

## GPT Cluster Analysis



# Естественные и сгенерированные



## Перспективно

The image features a background of numerous thin, wavy yellow lines that create a sense of motion and depth, resembling a liquid surface or a distorted grid. These lines are set against a solid black background. In the center of the image, there is a black rectangular box with a thin yellow border. Inside this box, the Russian word "Спасибо!" (Thank you!) is written in a bold, yellow, sans-serif font.

**Спасибо!**