

# Сравнительный анализ архивов научных журналов

Ф.В. Краснов (Газпромнефть НТЦ, 50+ научных статей, MBA, PMP, RHCE, к.т.н.),

М.Е. Шварцман (Библиотека имени Ленина, 90+ научных статей, к.т.н.),

А.В. Диментов (Национальный Электронно-Информационный Консорциум, [elpub.ru](http://elpub.ru))

# Суть

## **Бизнес потребность:**

Редакции журналов хотят понимать, что нужно изменить, чтобы войти в глобальные индексы цитирования (Scopus, WoS).

## **Наш подход:**

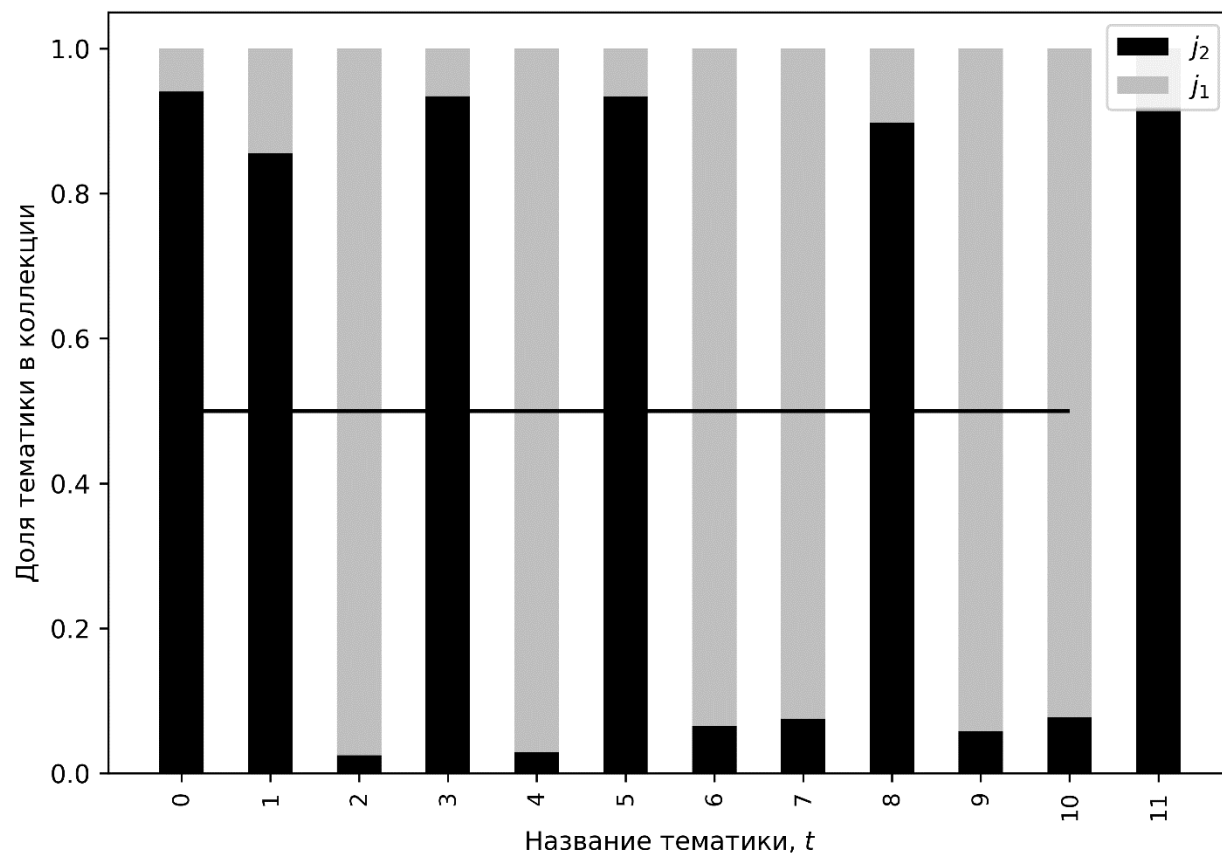
Бенчмаркинг архива журнала с мировыми лидерами с целью выявить отличия и выработать рекомендации по их устранению.

# Результат (1)

- Онлайн сервис по бенчмаркингу архивов журналов,
- Полностью автоматизированный процесс,
- Комплексная методика на основе анализа текстов и мета-информации о научных статьях.

***Сравнив архивы журналов, мы можем сказать редакции каких тематик, статей и авторов не хватает, чтобы удовлетворять требованиям Scopus и WoS.***

# Результат (2)



Коэффициент контентной аутентичности

Topic 0: hippocampal\_subfield, subiculum, subfield, hippocampal\_region, presubiculum, parasubiculum, hippocampal\_tail,

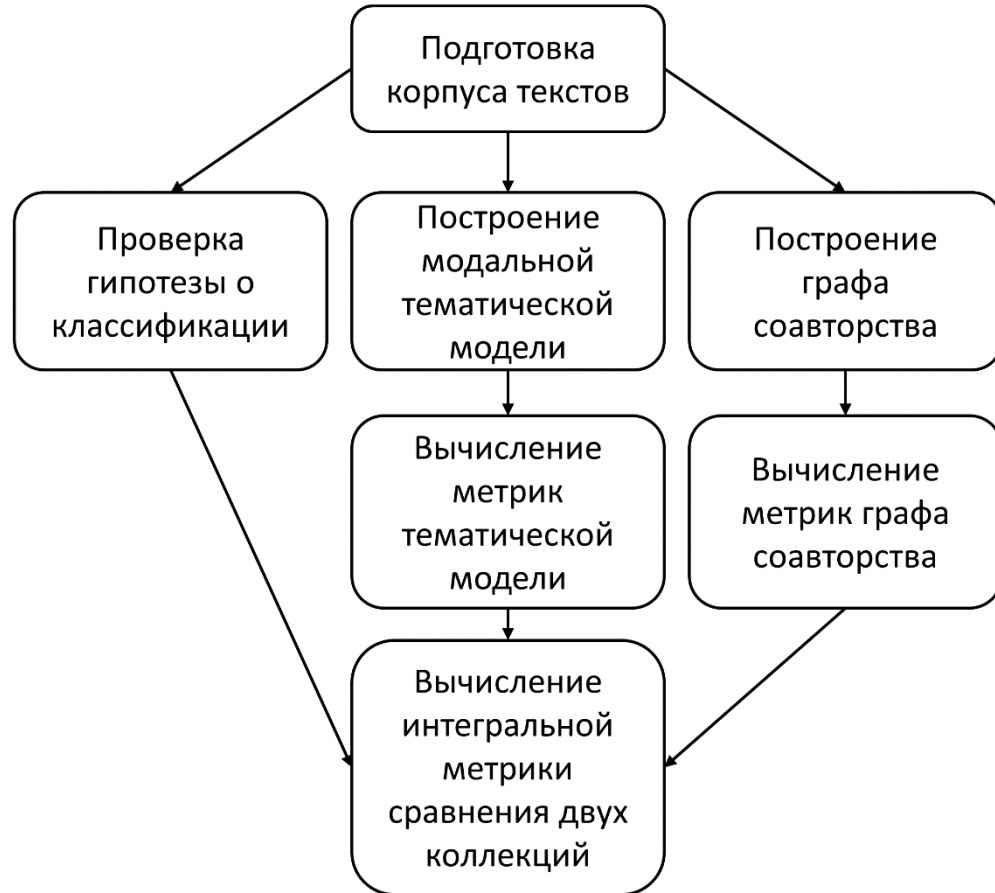
Topic 1: amygdala, fmri, ptsd, functional\_magnetic\_resonance\_imaging, activation, emotion\_regulation, ventromedial\_prefrontal\_cortex,

Topic 2: adhd, attention\_deficit\_hyperactivity\_disorder, adhd\_patient, intra\_individual\_variability, adolescent, erp, neurophysiological\_datum,

Topic 3: multiple\_sclerosis, flair, lesion, ms, segmentation, image, ms\_patient,

Topic 4: priori\_hypothesis, study\_design, publication, subject\_characteristic, medline, mr\_spectroscopy, magnetization\_transfer,

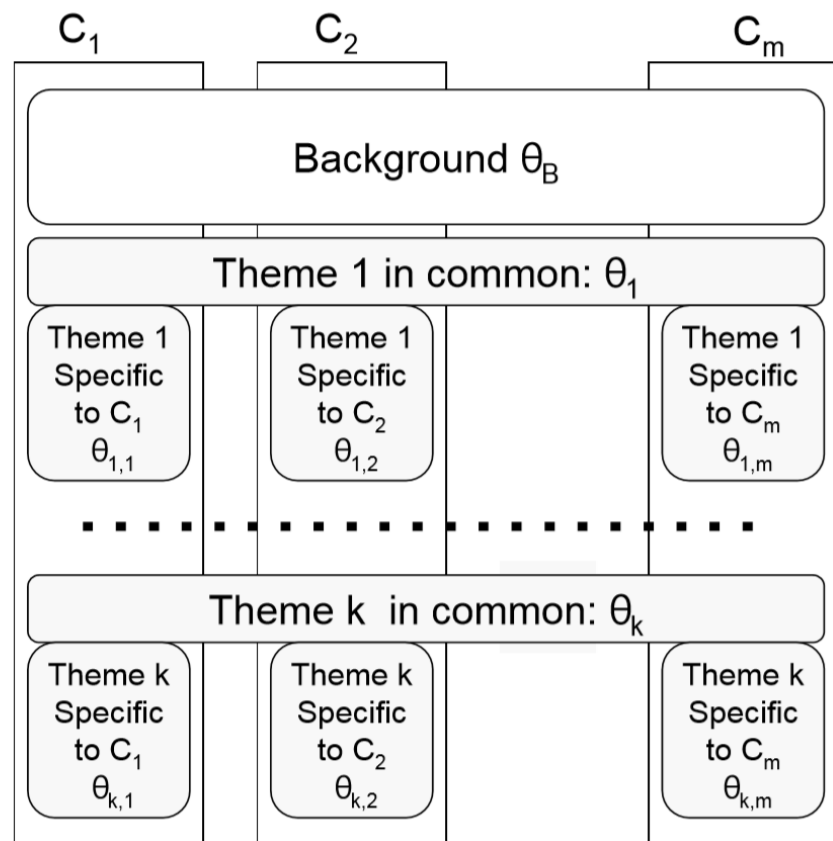
# Результат (3)



Автоматизированный процесс:

1. Сбор данных;
2. Гипер-параметры;
3. Уменьшение «шума»;
4. Уровень автоматизации.

# Методика Comparative Text Mining



ChengXiang ("Cheng") Zhai ([Short Biography](#))

Donald Biggar Willett Professor in Engineering

Department of Computer Science

(Also affiliated with [Carl R. Woese Institute for Genomic Biology](#), [Statistics](#), and [School of Information Sciences](#))

University of Illinois at Urbana-Champaign, 201 N. Goodwin Ave, Urbana, IL 61801

Office: 2116 Siebel Center ([directions](#))

Phone: [+1-217-244-4943](#) Fax: [+1-217-265-6494](#)

Email: [czhai AT illinois DOT edu](mailto:czhai@illinois.edu)

Assistant: Allison G. Mette ([agk AT illinois DOT edu](mailto:agk@illinois.edu), [+1-217-300-6504](tel:+1-217-300-6504), 2106 SC)

**Research Interests:** My general interests are in developing all kinds of novel **Intelligent Information Systems** (e.g., **intelligent search engines**, **recommender systems**, **text analysis engines**, and **intelligent task assistants**) to help people manage and exploit large amounts of data (i.e., "**big data**"), especially **text data**. I am particularly interested in building such intelligent systems for **improving health, medical care, education, and accelerating scientific discovery**, and building them based on theoretically sound frameworks, models, and algorithms that are also effective empirically. My **current interests** include *Intelligent Information Retrieval*, *Text Data Mining*, *Applied Machine Learning*, *Optimization of Human-Computer Collaboration*, *Biomedical and Health Informatics*, and *Intelligent Educational Systems*. For a **more detailed view** of my research interests, you may want to take a look at the slides of [some of my recent talks](#).

**Research Group:** [TIMAN](#), [DAIS](#)

**Publications & Recent Talks** ([Google Scholar Profile](#), [DBLP Profile](#))

The Book [Text Data Management: A Practical Introduction to Information Retrieval and Text Mining](#), ChengXiang Zhai and Sean Massung, ACM and Morgan & Claypool Publishers, July 2016, has now been translated into Chinese (see [the Chinese version](#)). The book is the main reference book for two **MOOCs on Coursera** that I'm teaching: (1) [Text Retrieval and Search Engines](#); (2) [Text Mining and Analytics](#).

Figure 2: The Cross-Collection Mixture Model

# Научная новизна

- Мягкая кластеризация с помощью тематической модели
- Тематическая модель для аннотаций (коротких текстов)
- Увеличение отношения сигнал/шум с помощью фокусировки на частях речи (POST, Part-of-speech tagging)

# Заключение

- Вывод на рынок нового аналитического продукта;
- От идеи до первой выручки 6 месяцев;
- Открытые программные каркасы (frameworks) можно продуктивно комбинировать, а не пытаться «допиливать»;
- Междисциплинарные команды экспертов высокоэффективны для инноваций.



Спасибо за внимание.

Федор Краснов: [krasnov.fv \(at\) gazpromneft-ntc.ru](mailto:krasnov.fv@gazpromneft-ntc.ru)

Михаил Шварцман: [shvar \(at\) rsl.ru](mailto:shvar@rsl.ru)

Александр Диментов: [dimentov \(at\) neicon.ru](mailto:dimentov@neicon.ru)