

Вычислительный Центр Российской Академии Наук

**«Выявление и визуализация
метрических структур на множествах
пользователей и ресурсов Интернет»**



Постановка задачи

- **Разработать методы выявления структур сходства в клиентских средах**
- **Специфицировать методы для обработки логов поисковой машины**
- **Построить карты сходства ресурсов Интернет**



Анализ Клиентских Сред

Задачи:

- **принятие решений**
- **маркетинг**
- **управление взаимоотношениями с клиентами**



Анализ Клиентских Сред

Схема:

- **Исходный протокол действий клиентов**
- **Выполнение цепочки процедур**
- **Решение конкретной задачи**

Исходные данные - лог поисковой системы

```
100002171080304956
  Валютное регулирование и валютный контроль в РФ 1110371715 281743 0
    http://tandem-forum.ru/seminars/info/seminar39.html 1110371732
    http://tandem-law.ru/seminars/info/seminar283.html 1110371757
    http://www.nns.ru/krizis/kabinet/2104/krizis8.html
    http://dit.perm.ru/articles/management/data/021217.htm 1110371805
    http://www.pomosch.com/article.php?sectionId=6&articleId=53
    http://www.e-ducate.ru/info/3523
    http://www.seminarium.ru/cgi-bin/viewer/events/data.cgi?eventID=1007519
    http://anons.webplan.ru/cgi-bin/shownew.cgi?numn=1518
    http://www.kadis.ru/daily/index.html?id=4405
    http://www.e-seminar.ru/root.phtml?PageID=168&EventID=486&CenterID=165 1110373262
  работа 1110527366 97008639 0
    http://www.job.ru/
    http://www.joblist.ru/
    http://rabota.mail.ru/
    http://www.zarplata.ru/ 1110527399
    http://www.ustroim.ru/
    http://www.rabota.ru/
    http://www.infocentre.net.ru/
    http://www.alljobs.ru/
    http://job.kulichki.net/
    http://www.webit.ru/
100002701110038640
  ...
```

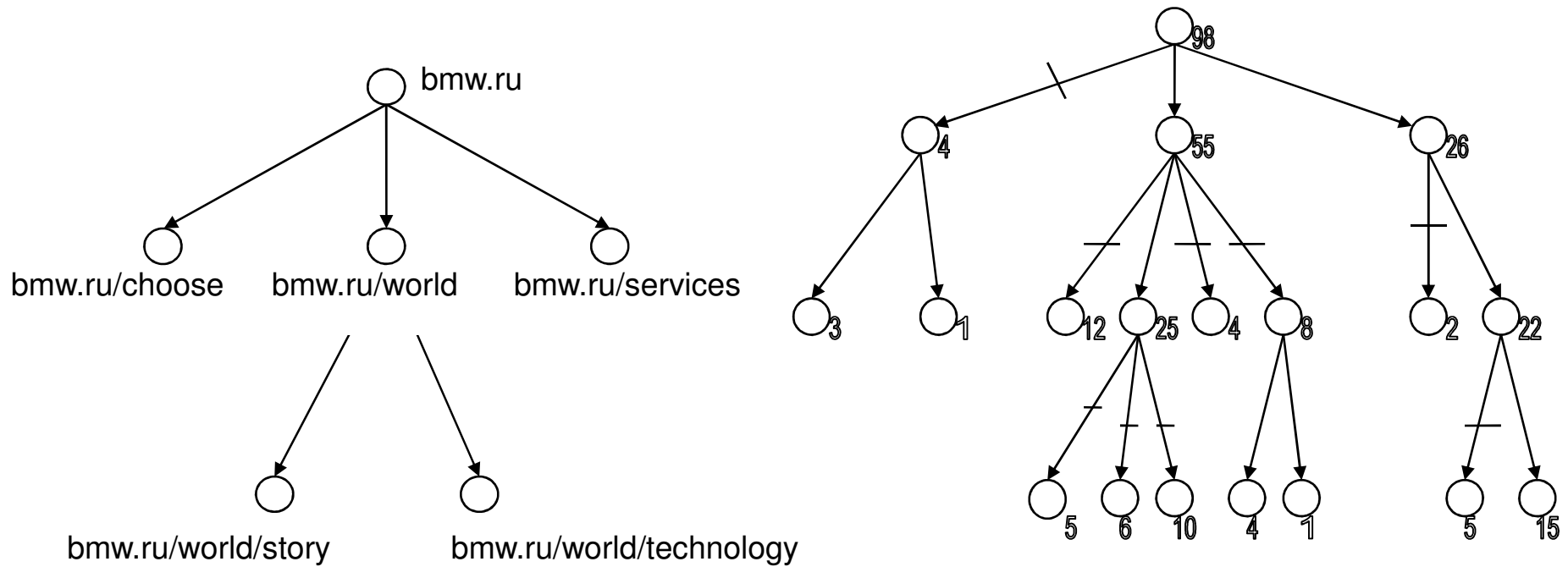
- 14 606 пользователей
- 1 972 636 ресурсов (из них 129 600 были выбраны пользователями)
- Интервал времени: 1 неделя работы поисковой системы
- Объем лога 4Гб (время обработки около 30 секунд)



Общая схема обработки данных

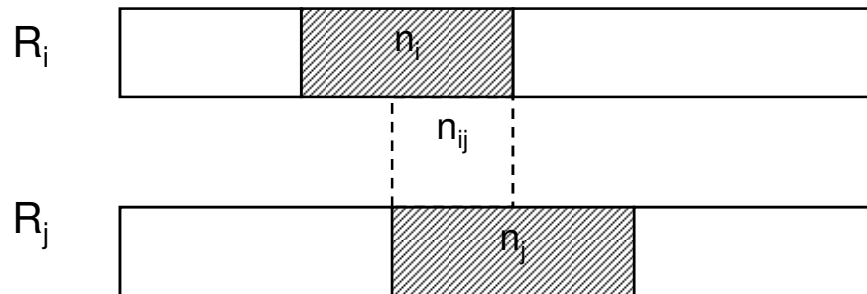
1. Формирование полных словарей пользователей и ресурсов
2. Редукция словарей
3. Формирование матрицы посещений
4. Формирование разреженной матрицы сходства ресурсов
5. Построение карты сходства всех ресурсов (карта Интернета)
6. Построение локальной карты сходства в окрестности заданного ресурса.

Формирование и редукция словарей пользователей и ресурсов

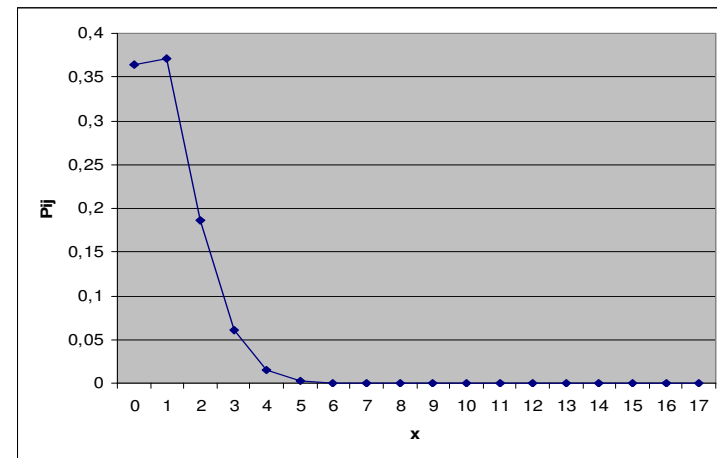


Пример построения и редукции словаря ресурсов с порогом 13 посещений

Вычисление оценок сходства ресурсов



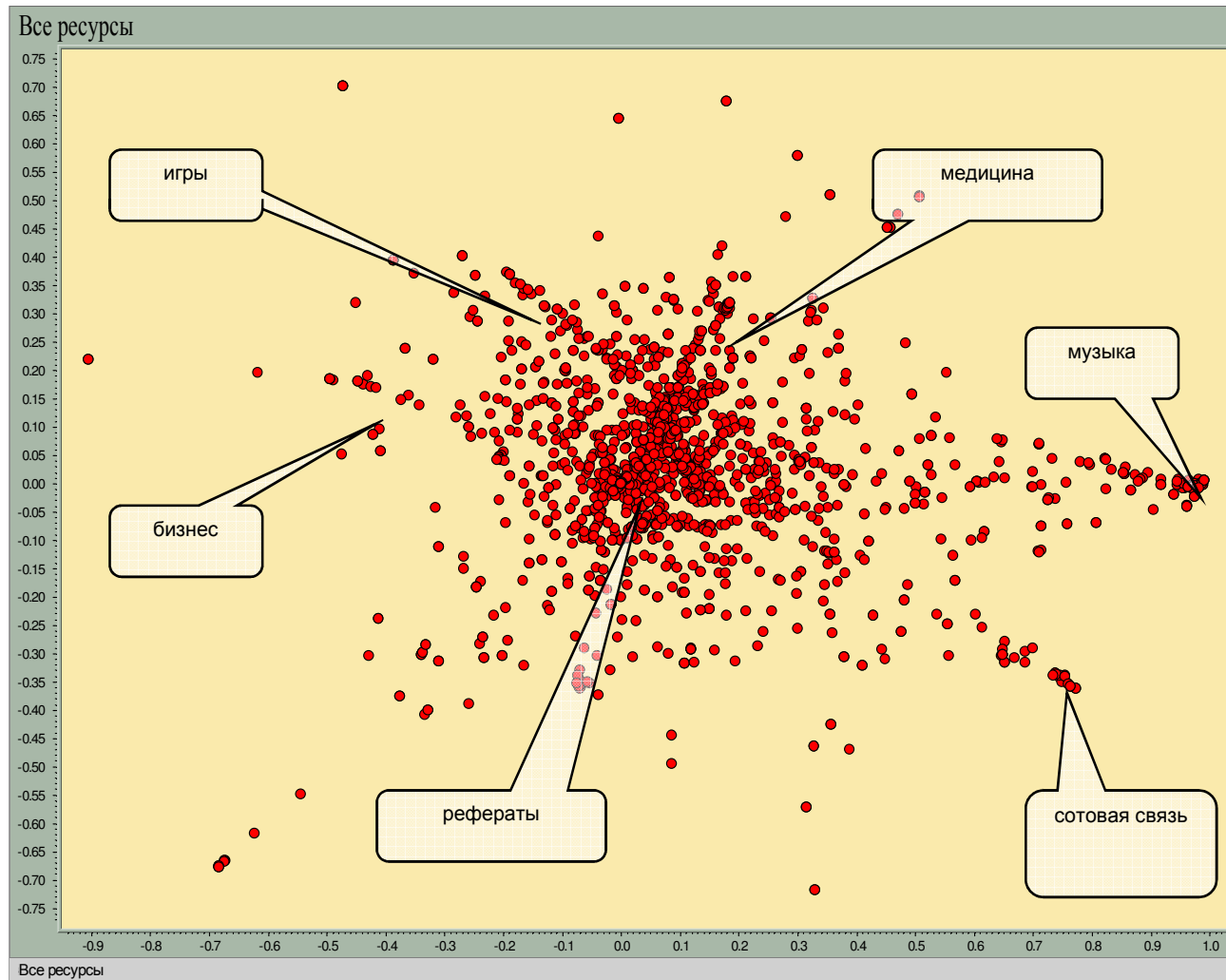
Гипергеометрическое распределение:



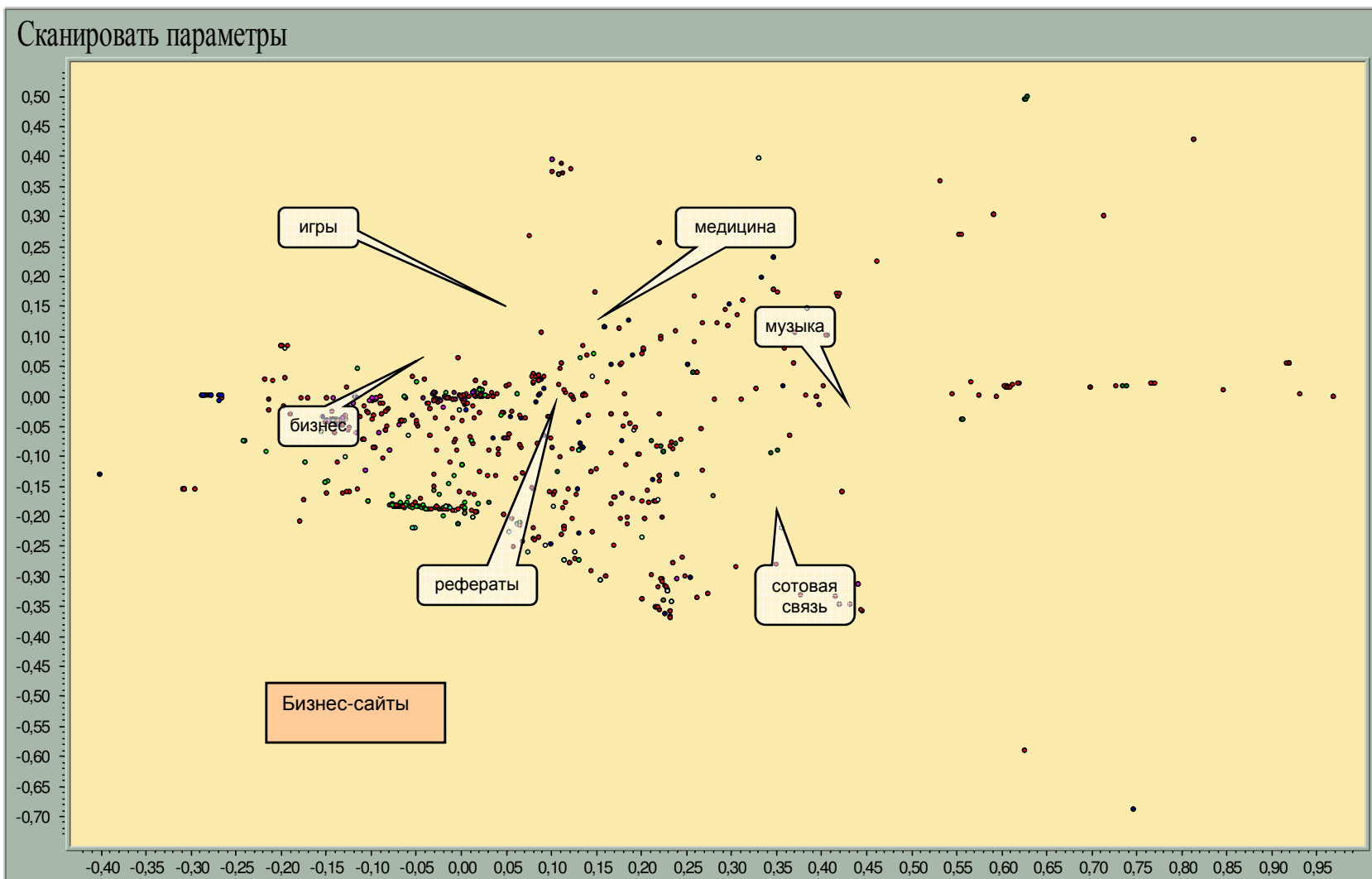
$$P_{ij} = P(n_{ij} = x) = \frac{C_{n_i}^x C_{U-n_i}^{n_j-x}}{C_U^{n_j}}$$

$$\rho(i, j) = \left(\frac{|\ln \alpha|}{|\ln P_{ij}|} \right)^3$$

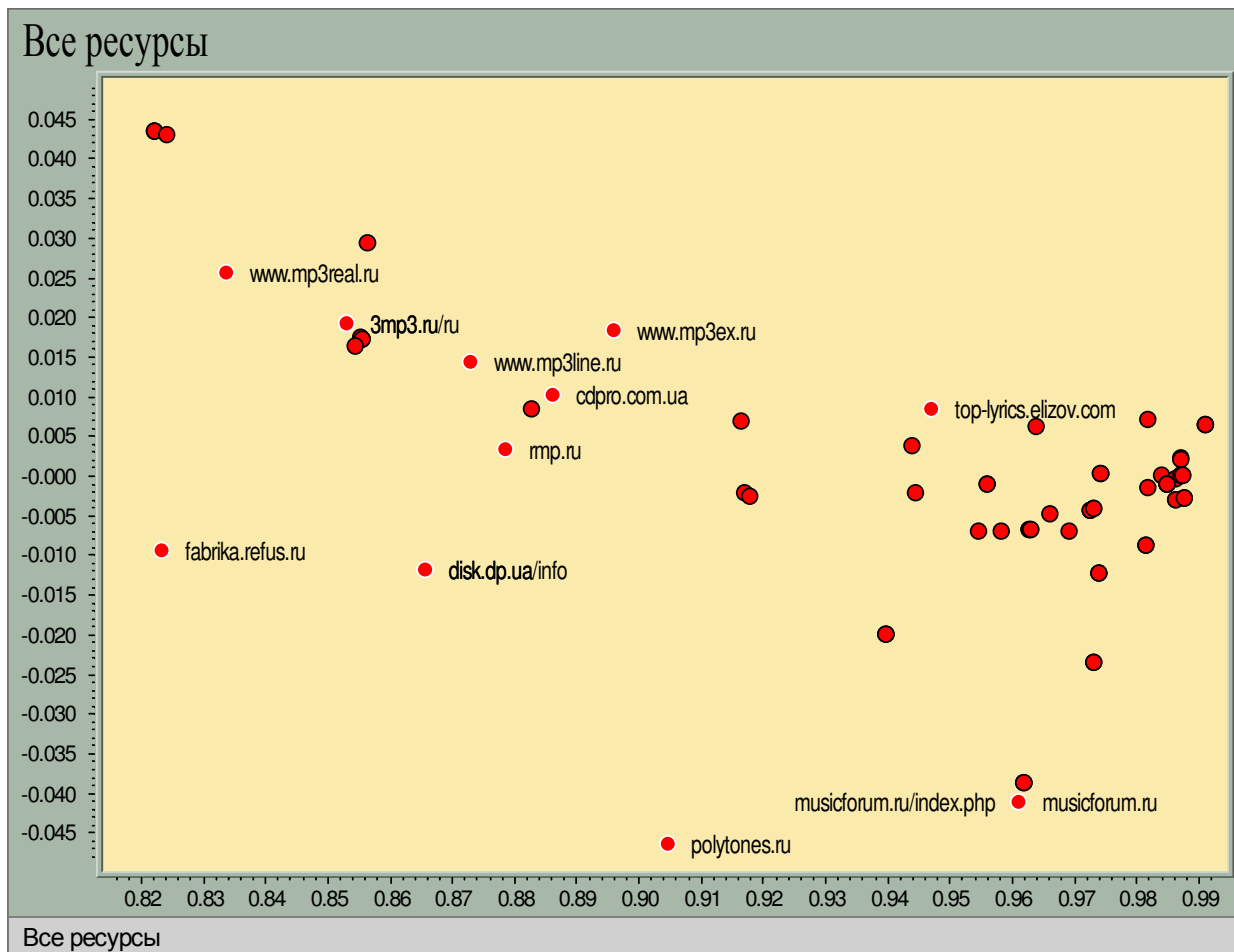
Карта сходства ресурсов



Карта сходства после оптимизации параметров

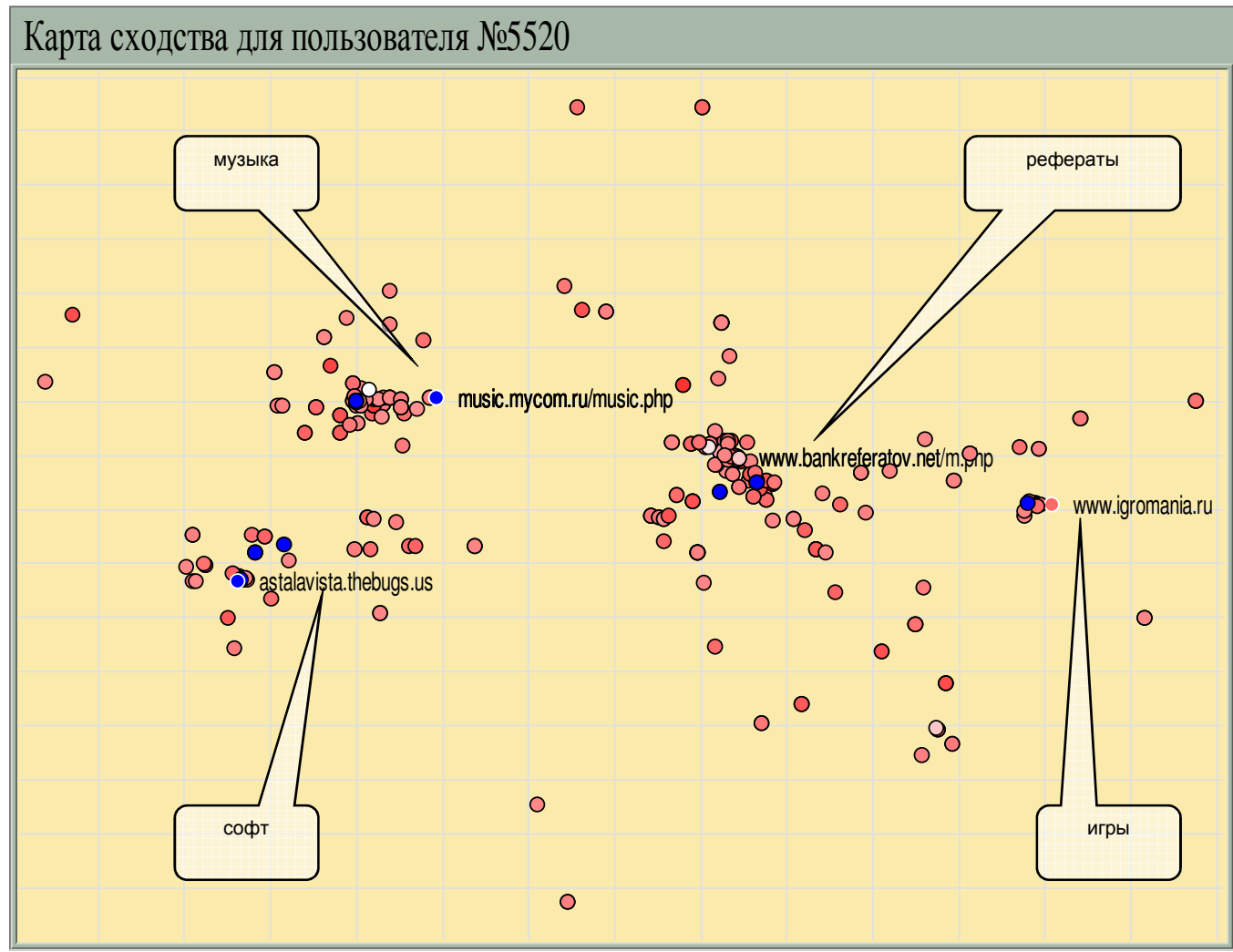


Фрагмент карты сходства



- | | |
|--|--|
| my3blka.narod.ru | www.mp3zzz.ru |
| mp3load.ru | mp3-unique.liveinternet.ru |
| mp3forum.ru | www.mp3collection.ru |
| mp3rank.ru | www.mp3real.ru |
| top-lyrics.elizov.com | rmp.ru |
| www.mp3ex.ru | mp3.allabout.ru |
| www.mp3line.ru | www.filesale.ru |
| www.oxid.ru | www.agharta.ru |
| 3mp3.ru | moremusic.com.ru |
| www.miditext.ru | ru.mp3s-download.com |
| www.rusmuz.ru | meloman.tehnofil.ru |
| www.mp3search.ru | www.mp3tones.ru |
| thesearch.ru | www.muZZone.com |
| www.mp3.myfind.ru | www.mp3zone.ru |
| www.russianmusic.de | www.zaycev.net |
| www.mp3host.ru | disk.dp.ua |
| www.mp3up.ru | www.uliss.ru |
| mp3.volpi.ru | musicforum.ru |
| cdpro.com.ua | |

Персональная карта сходства





Основные результаты

- Разработаны и реализованы эффективные алгоритмы анализа логов поисковой машины
- Построена карта сходства ресурсов Интернет
- Предложена методика построения локальных карт сходства