



The Crawler

Alexey Voropaev

go.mail.ru



- <http://go.mail.ru>

- Web search for:

- Russian
- Ukrainian
- Kazakh

- 9% of market share

- 100+ developers



хью лори фото

Интернет Картинки Видео Новости Обсуждения Ответы

Картинки



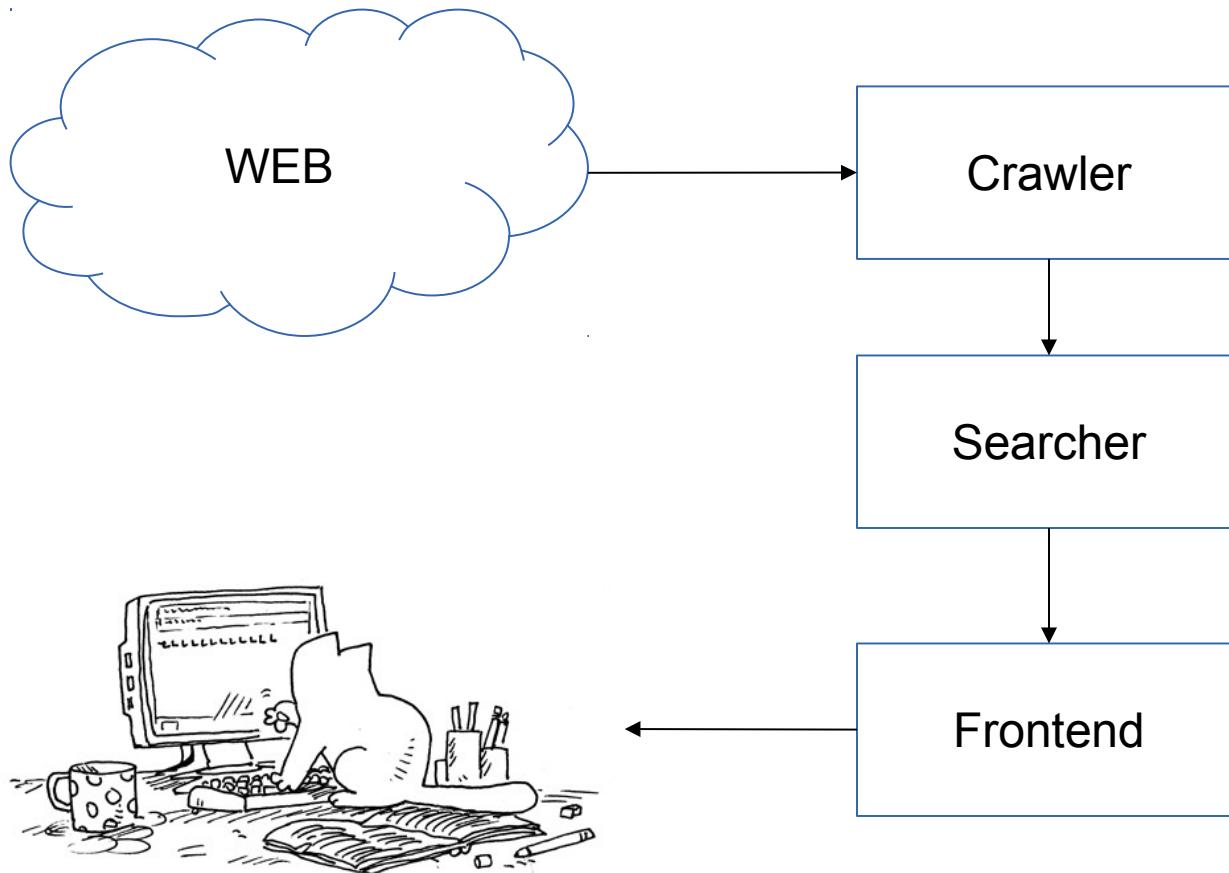
[Лори, Хью — Википедия](#)
ru.wikipedia.org/wiki/Лори,_Хью

Родители Хью Лори — шотландцы. ... Хью Лори и актриса Имельда Стонтон дважды появлялись на экране в роли мужа и жены: в фильме «Разум и чувства» (1995) и «Друзья Питера» (1992). ... Фотографии разных лет.



[Хью Лори](#)
afisha.mail.ru
Фотографии

Актер Хью Лори родился в семье врача и домохозяйки. Учился в частных школах, окончил Кембриджский университет. В университете был участником любительского театра Footlights Dramatic Club. Уже в...



«Running a web crawler is a challenging task»

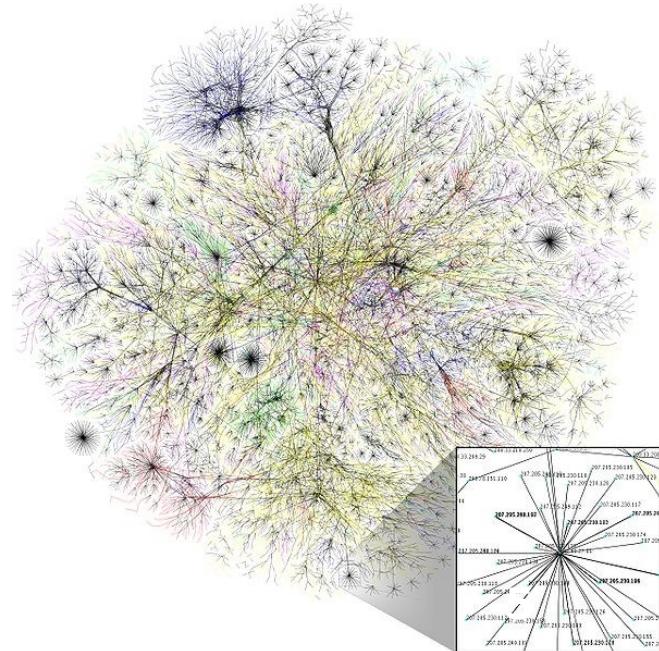
Sergey Brin and Lawrence Page, 1998

The web:

- Web sites: $\sim 10^7$
- Web pages (sensible): $\sim 10^{10}$
- URLs: ∞
- 60% of pages are located on 10000 biget sites

Search engine:

- Limited index: $\sim 10^9$ pages
- Limited bandwidth: max ~ 100 pps
- Limited ranking: $\sim 10^5$ ppq



Scheduler (aka long-term scheduler)

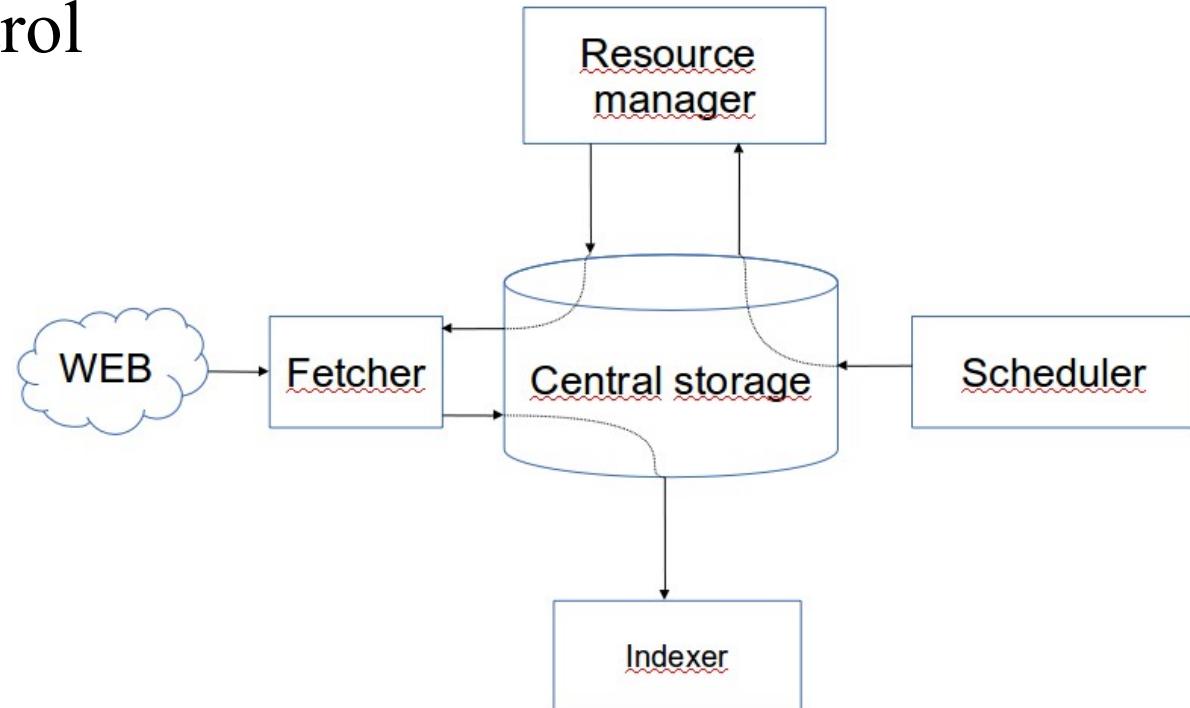
- Selection policy
- Refresh policy

Resource manager (aka short-term scheduler)

- Bandwidth control

Fetcher

- Just WGET



Indexer

- Analyze of downloaded content

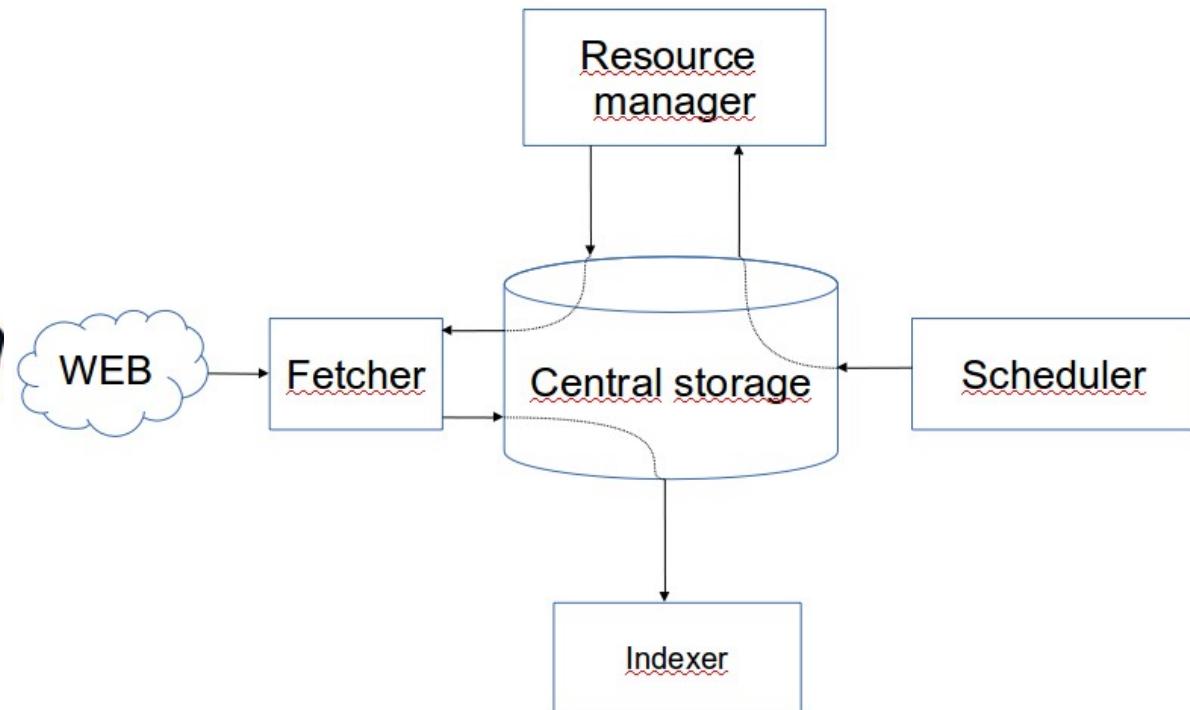
Auxiliary

- Logging
- Statistic

Central storage



APACHE
HBASE



Actually it is set of schedulers

- Web search
 - Download scheduler
 - Indexing scheduler
 - Discovery scheduler
 - Sandbox scheduler
 - Analyzer scheduler
- Image search
 - HTML scheduler
 - Image scheduler
- Experiments

Focused crawler – download pages with given topic.

We have crawling framework:

1. Define what is *interesting for crawling* pages
2. Find *initial seed*
3. Exploit *topical locality* to find other pages

What is “interesting for crawling” for Web-search?

- Page with Qlinks is interesting
- We can say nothing about page without Qlinks

– Link based

if A->B then B similar to A

– URL based

^http://aldebaran.ru/lov/[a-z]+/[a-z]+[0-9]+/\$

*

^http://aldebaran.ru/kid/krapiv/krapiv[0-9]*\$

*

^http://materinstvo.ru/\$
+id+module
~module=articles

– Content based

title="Анкета заблокирована"

query="Что вы думаете об этом товаре?"

For given:

- α fraction of urls from some group on a site
 N sampling size

Probability to find less then k urls:

$$p_{N,k}(\alpha) = \sum_{i=1}^k \binom{i}{N} \alpha^i (1 - \alpha)^{(N-i)}$$

$$P_{1000,10}(0.01) \approx 0.58$$

$$P_{1000,10}(0.02) \approx 0.01$$

$$P_{1000,10}(0.03) \approx 2 \times 10^{-5}$$

1. Sample N urls randomly
2. Generate features for each url
 - number of segments
 - list of parameters
 - <parameters=value>
 -
3. Select features with frequency αN
4. Do clustering:
 - Jaccard distance measure $K(a, b) = \frac{|A \cap B|}{|A \cup B|}$
 - Stack clustering

1. **^/wiki/File:[^/]+\.jpg\$**

/wiki/File:Spongilla_lacustris.jpg

2. **^/wiki/[^/]+\.jpg\$**

/wiki/Image:Deve.jpg

3. **^/wiki/Category:[^/]+\$**

/wiki/Category:Roman-era_historians

4. **^/wiki/Talk:[^/]+\$**

/wiki/Talk:North_Light

...

- 404, 500, timeout etc.
- spam, porno etc.
- duplicate detection
- news
- navigation removal
- etc.

LENTA.RU Россия
7 ФЕВРАЛЯ 2013, ЧЕТВЕРГ, 23:30

Все Политика Общество Преступность Дороги Происшествия

19:35, 7 февраля 2013

Москва глухих

Что не так с сотовыми операторами в Москве

Уникальный контент



LENTA.RU Культура
7 ФЕВРАЛЯ 2013, ЧЕТВЕРГ, 23:30

Все Кино Музыка Театр Книги Архитектура Искусство Игры

19:07, 7 февраля 2013

Шотландия превыше всего

Шон Коннери поведал в автобиографии о своей главной страсти



ПОСЛЕДНИЕ НОВОСТИ

20:47 Путину доложили о части билетов на Олимпиаду

22:54 Сборная Норвегии пришла первой в соревнование по биатлону

21:34 Директора фабрики вызвали на допрос

ПОСЛЕДНИЕ НОВОСТИ

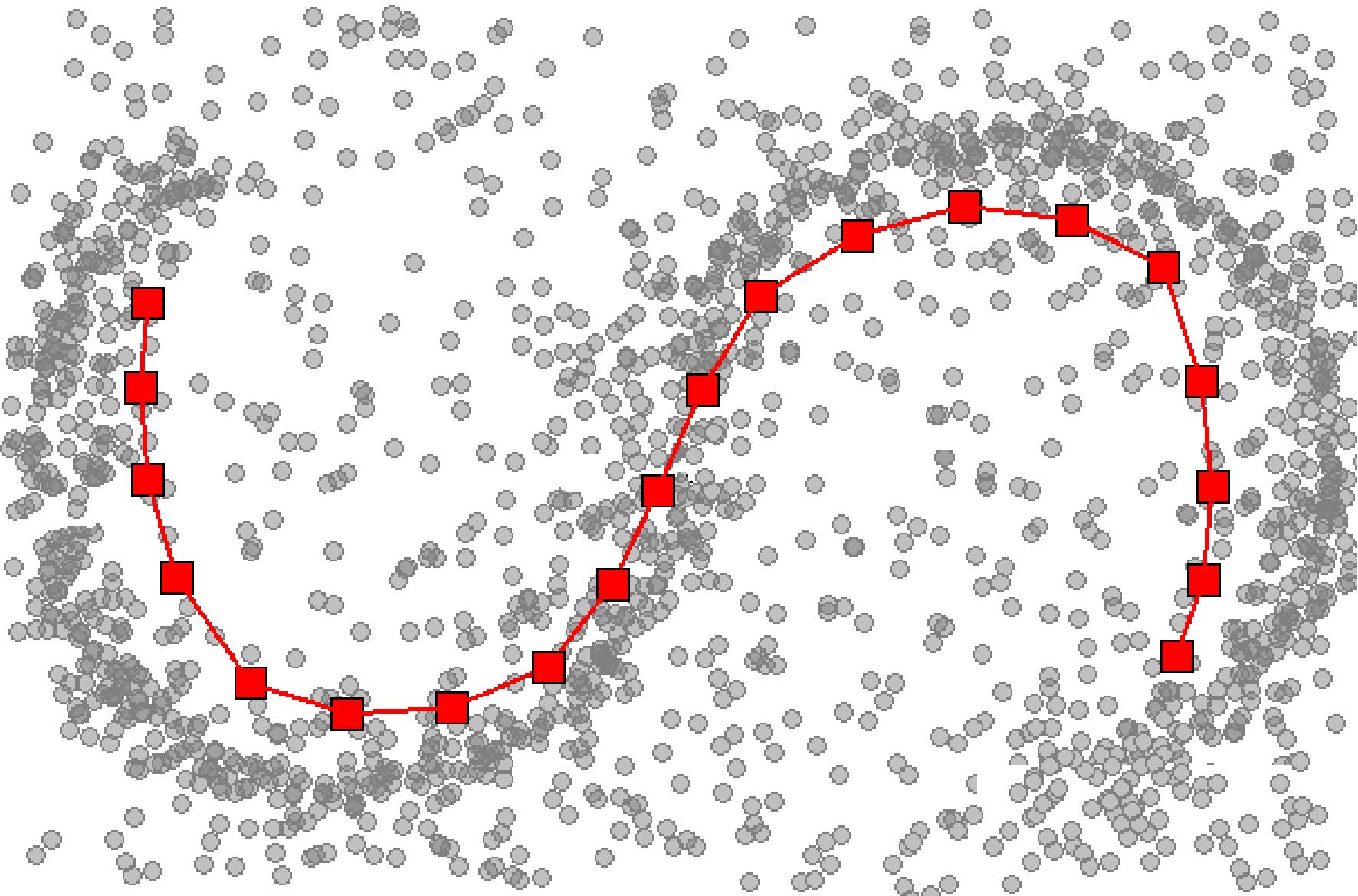
20:47 Путину доложили о части билетов на Олимпиаду

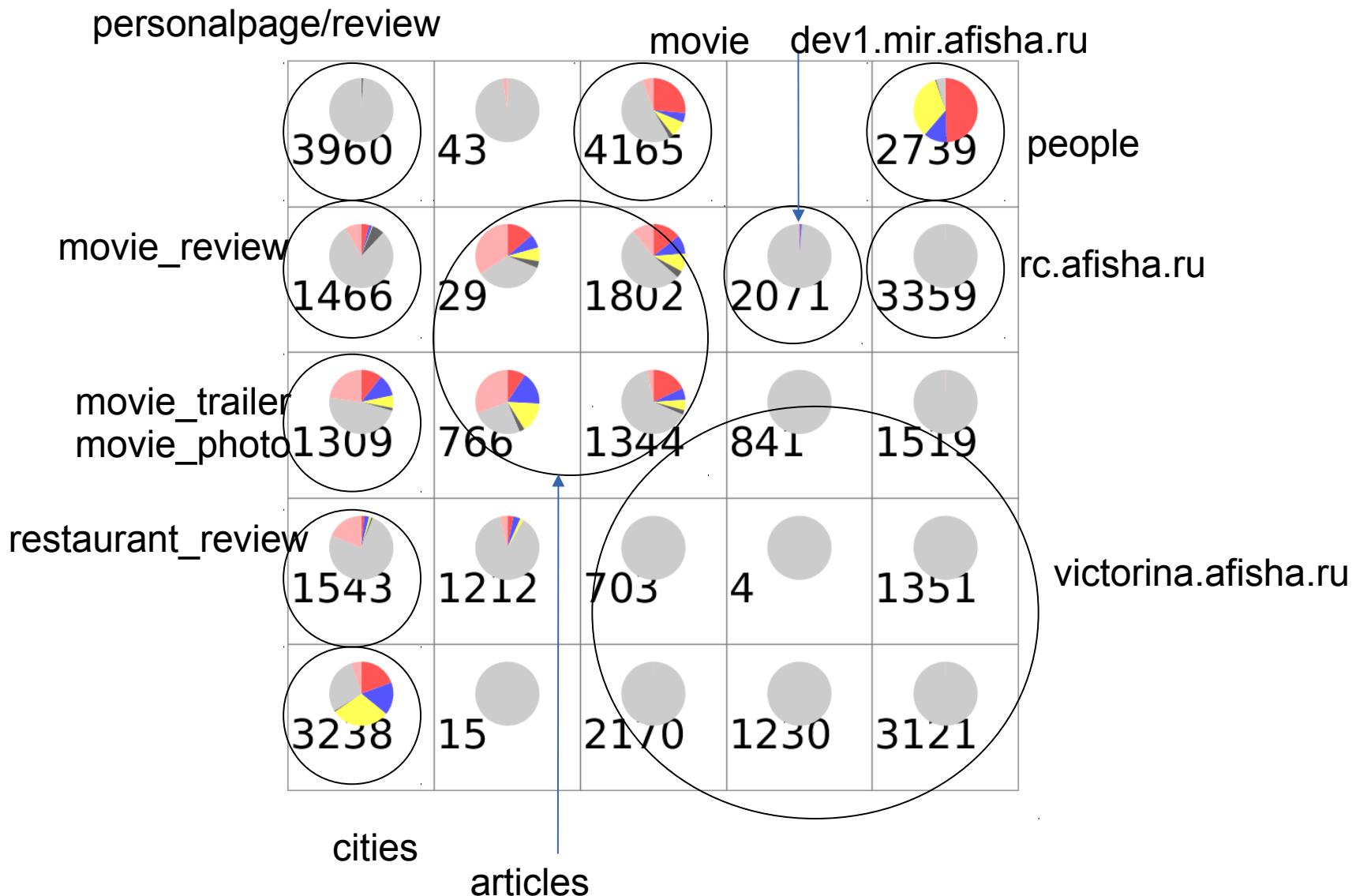
22:54 Сборная Норвегии пришла первой в соревнование по биатлону

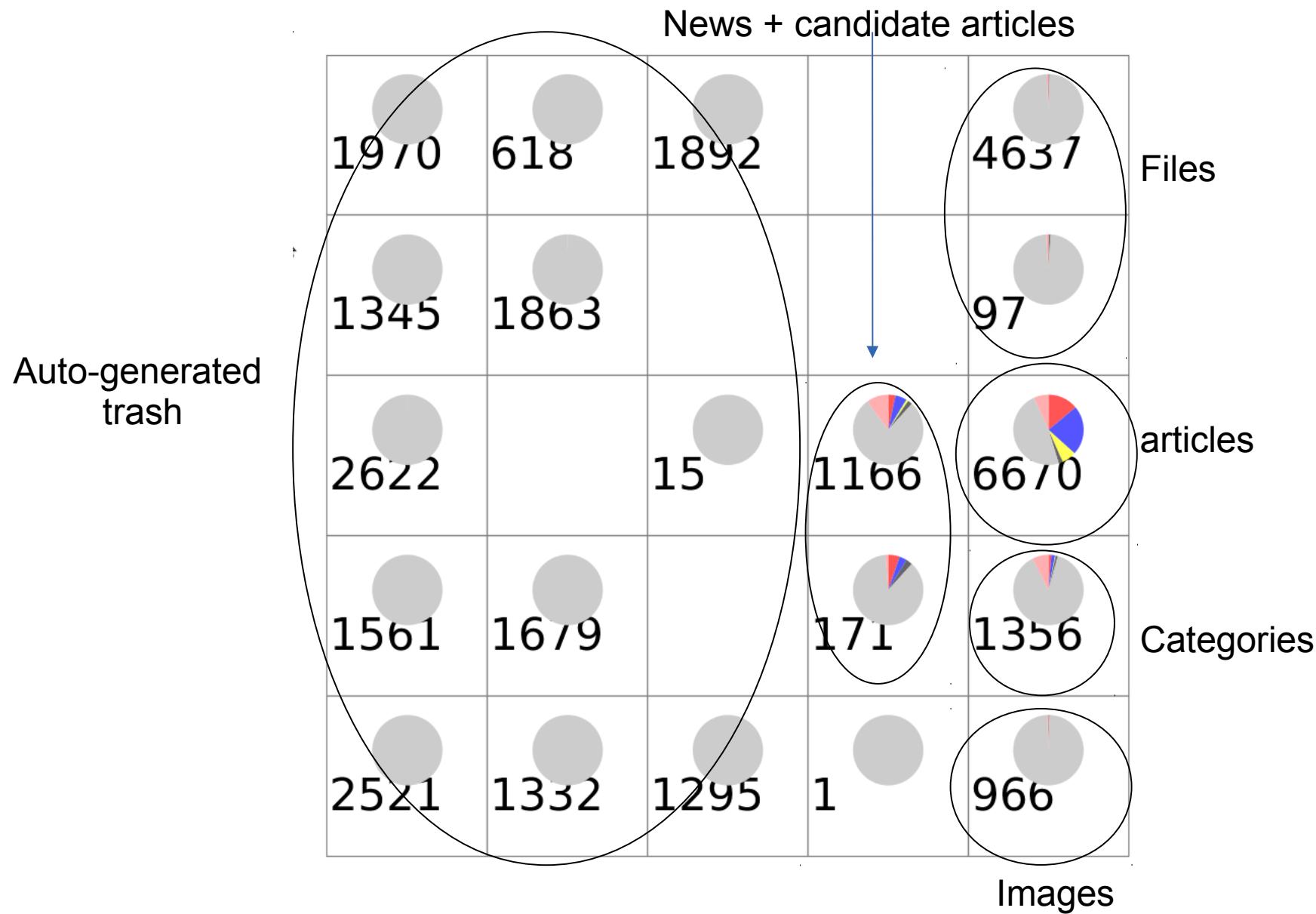
21:34 Директора фабрики вызвали на допрос

SOM / Kohonen map

go.mail.ru







Goal: to build index with fixed size

Site:

- for hosting (about 1000 sites in ru) domain level 3
- otherwise domain level 2

Quota:

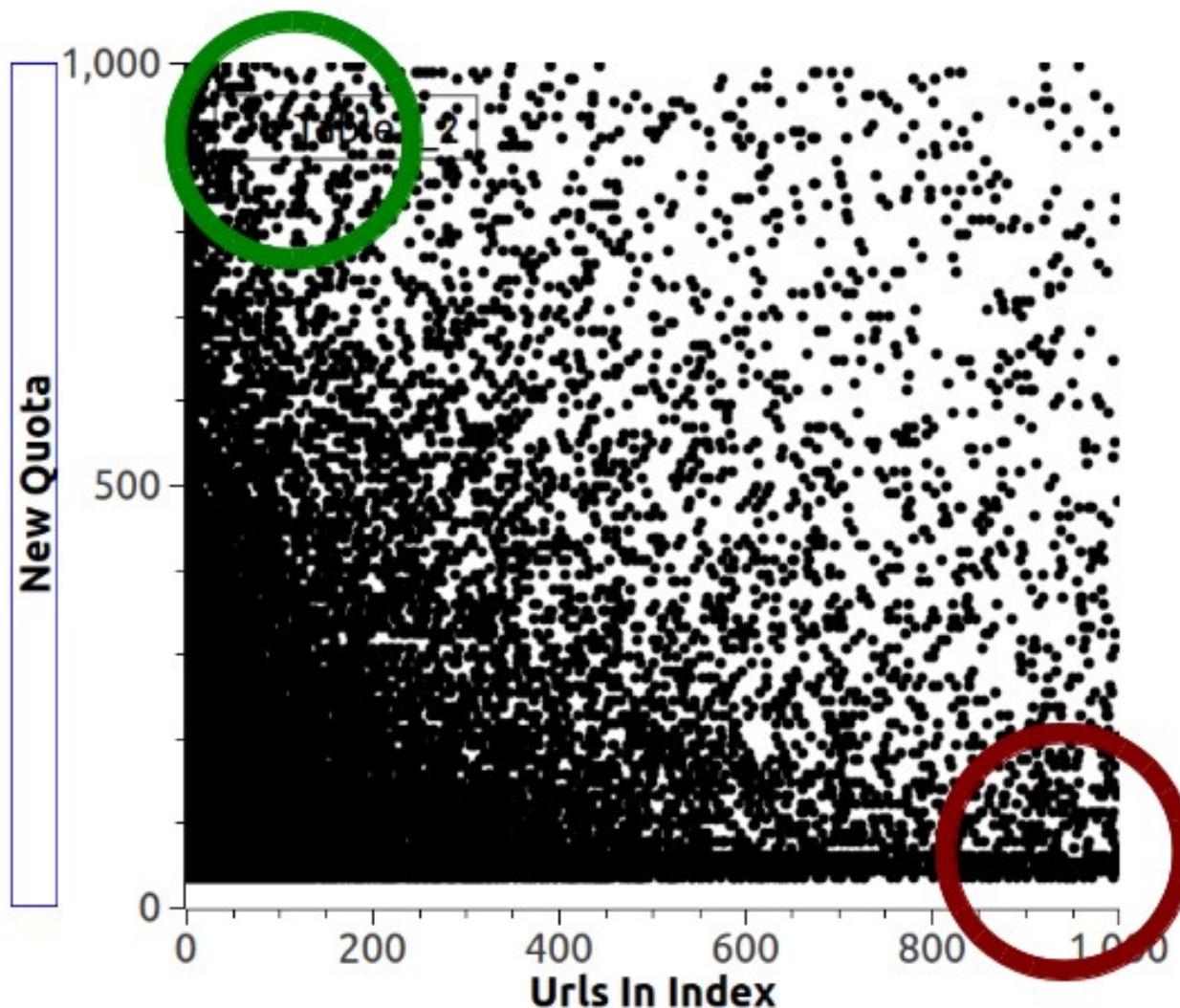
- MIN_QUOTA ~ 100
- QUOTA = #Pages with Qlinks * MIN_QUOTA

Diversity:

- We do quoting between stones inside of big sites

Quoting: compare with old algorithms

go.mail.ru



Good
Blogs

blogspot/livejornal

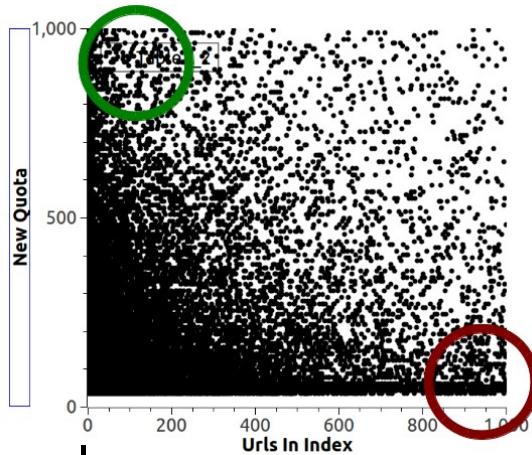
Good sites

use4blog.com

gagadget.com

Popular foreign sites

Robots, ban, etc.



Bad
Trash sites

Скачать сочинение на тему как повысить речевую культуру - Вы искали:

- 1. [как перепланировать двухкомнатную хрущевку](#)
- 2. [Агабекян И рабочая программа для спо](#)
- 3. [образ святослава](#)
- 4. [видео эротичные девушки в чулочках](#)
- 5. [ключ для 16-фло Кеттела](#)
- 6. [windows 7 professional cis and ge](#)
- 7. [Бесплатный конвертер pdf в excel](#)
- 8. [решебник грамматика английского языка 5 класс барашкова](#)
- 9. [спесня царевны забавы» ноты фортепиано](#)
- 10. [сочинение в драгунский](#)
- 11. [загадки про животных](#)

Features:

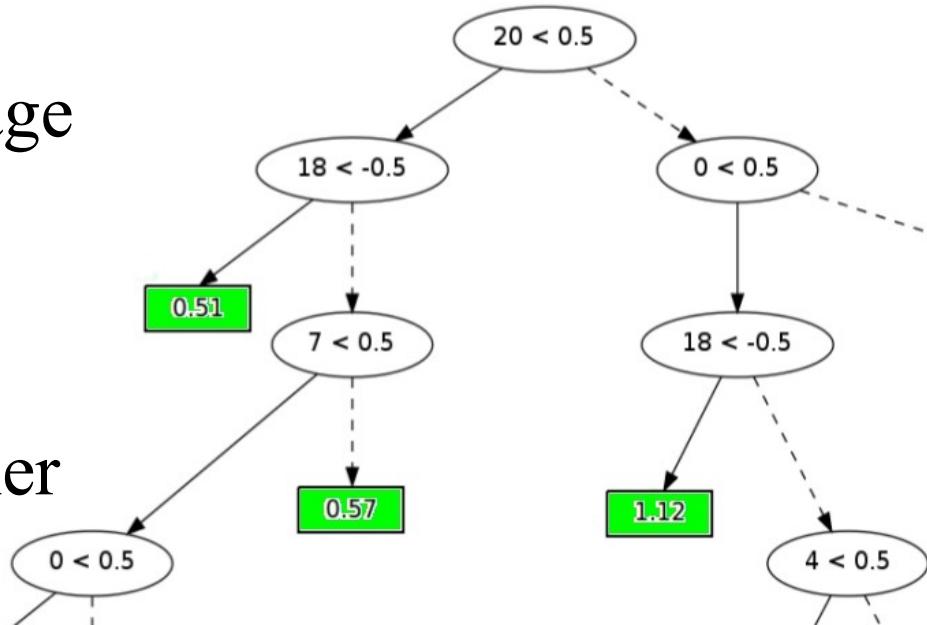
- Sekitei features
- Link features (Indegree, PR, etc.)
- Antispam features (e.g. #links from bad sites)

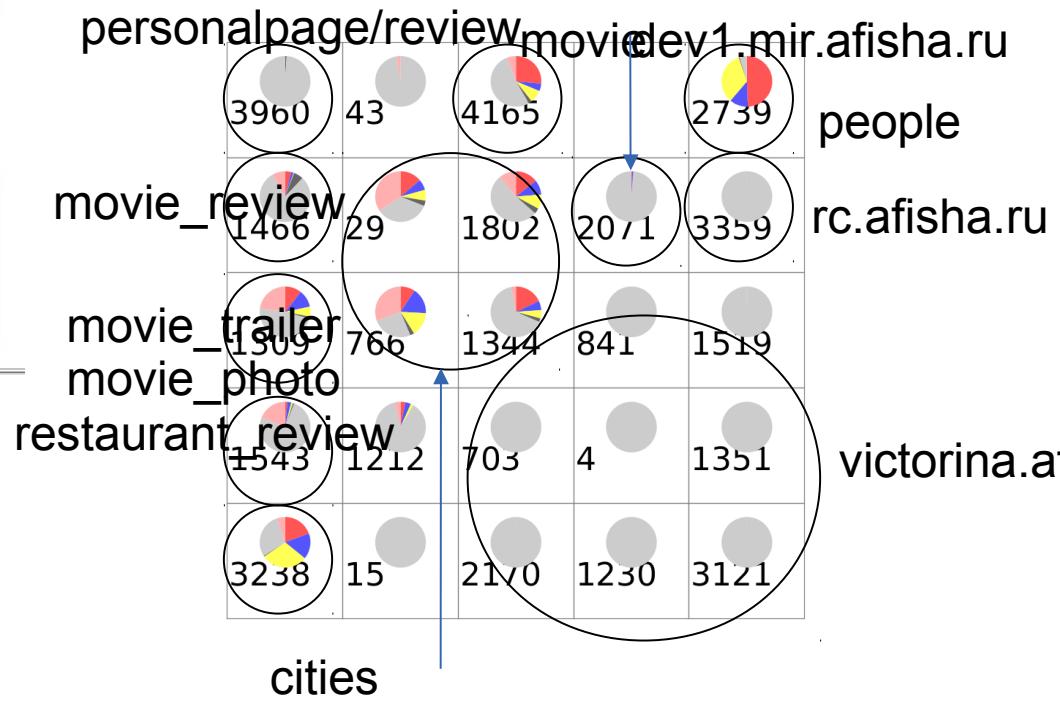
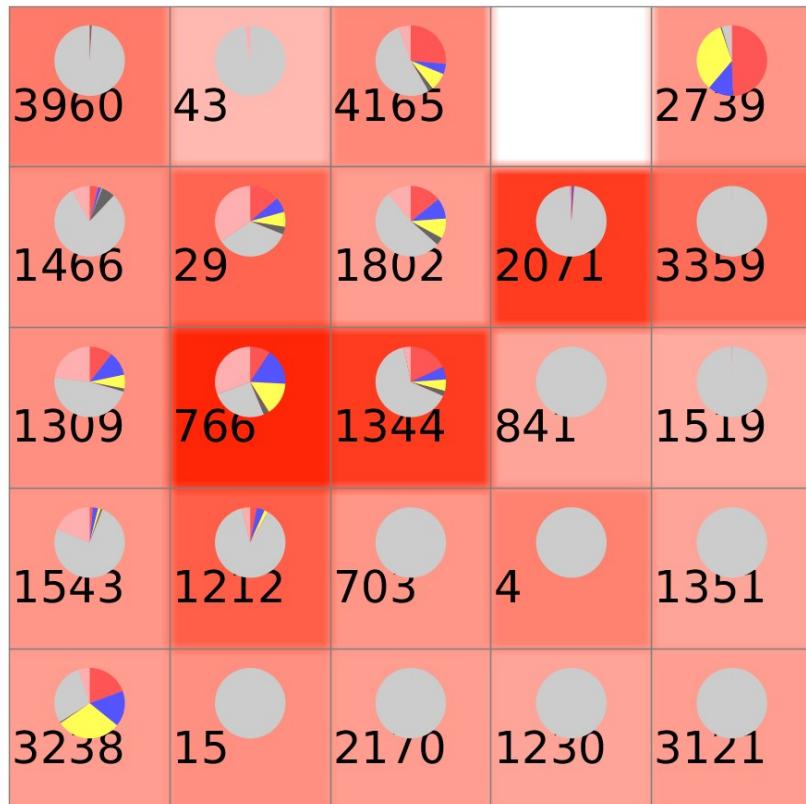
Model: gradient boosting decision trees

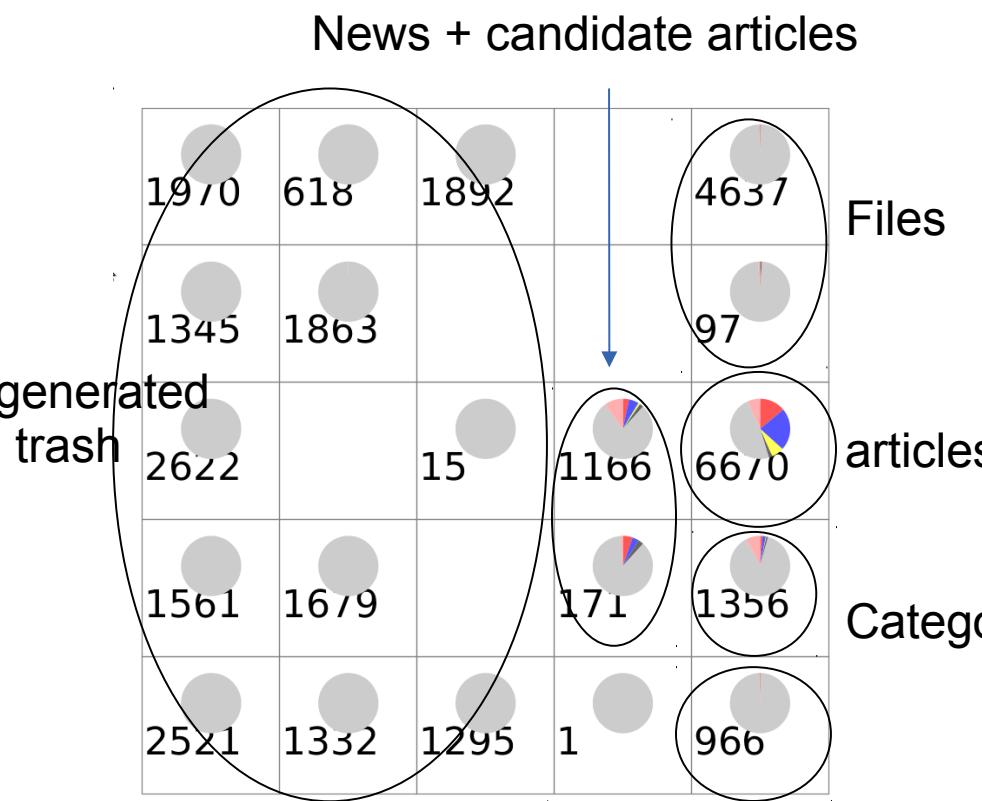
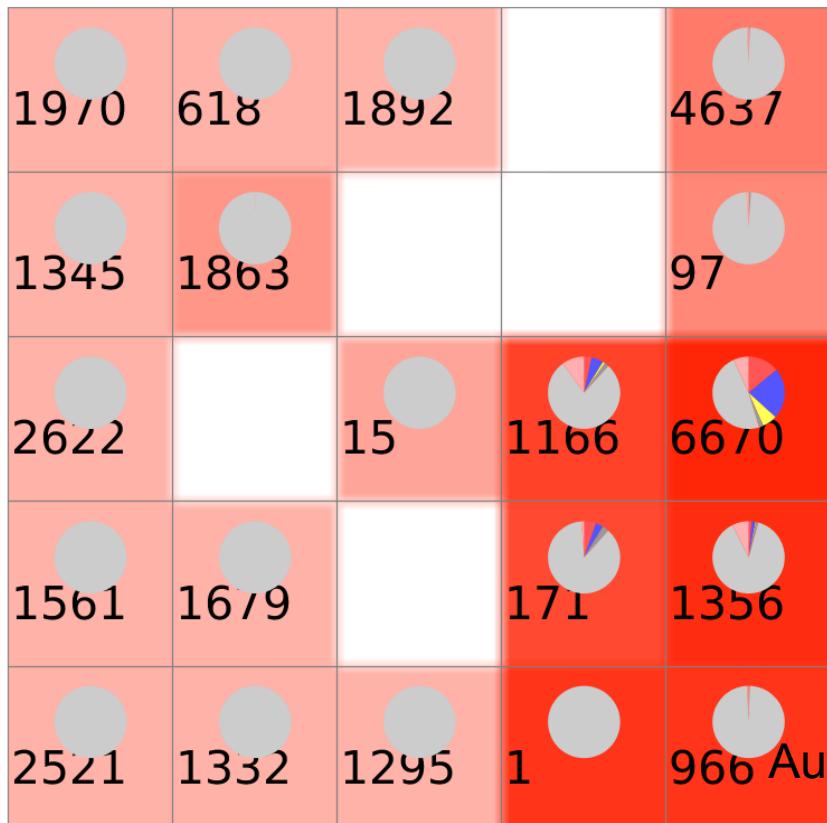
Target: predict #qlink for a page

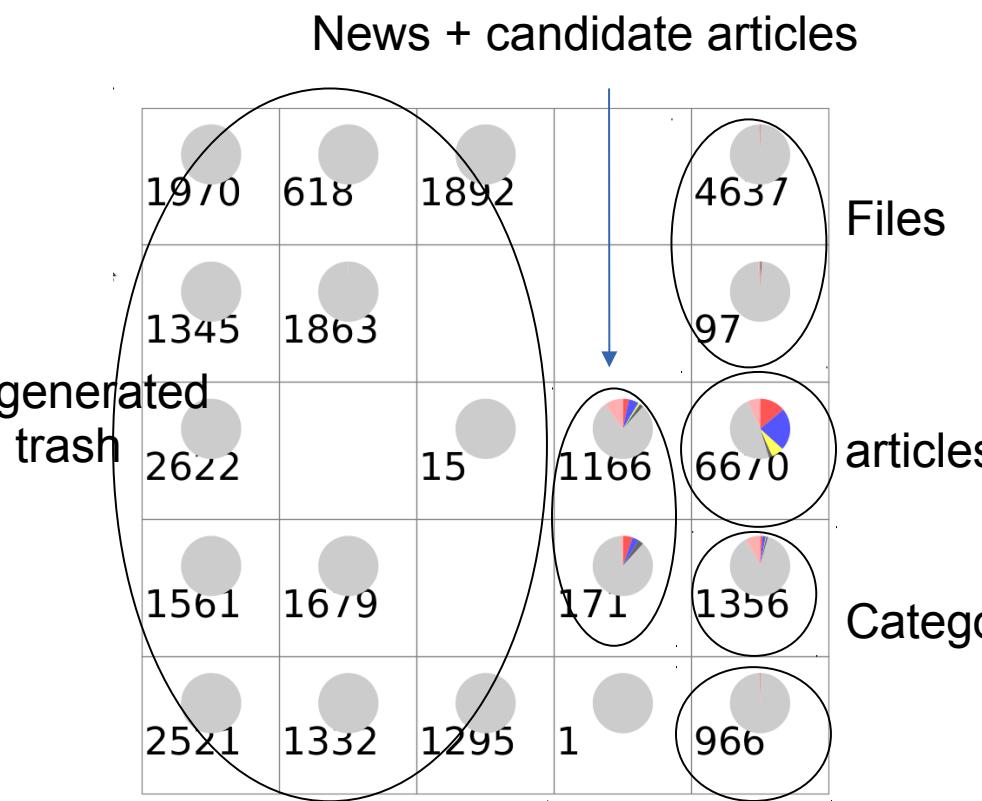
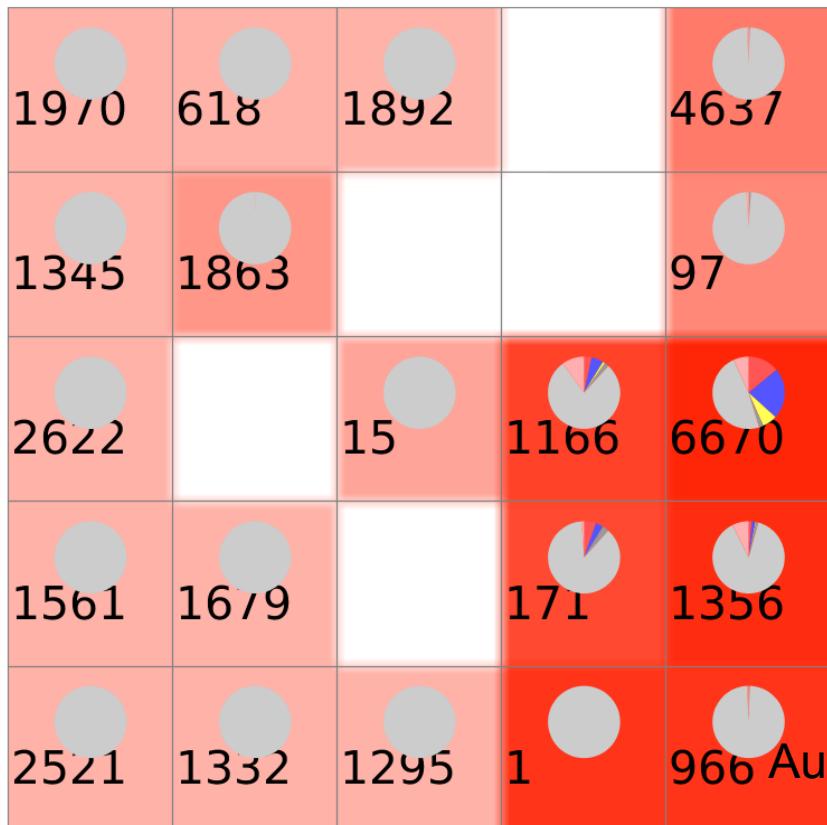
2 kind of models:

- individual (for big sites)
- one “wide” model for other



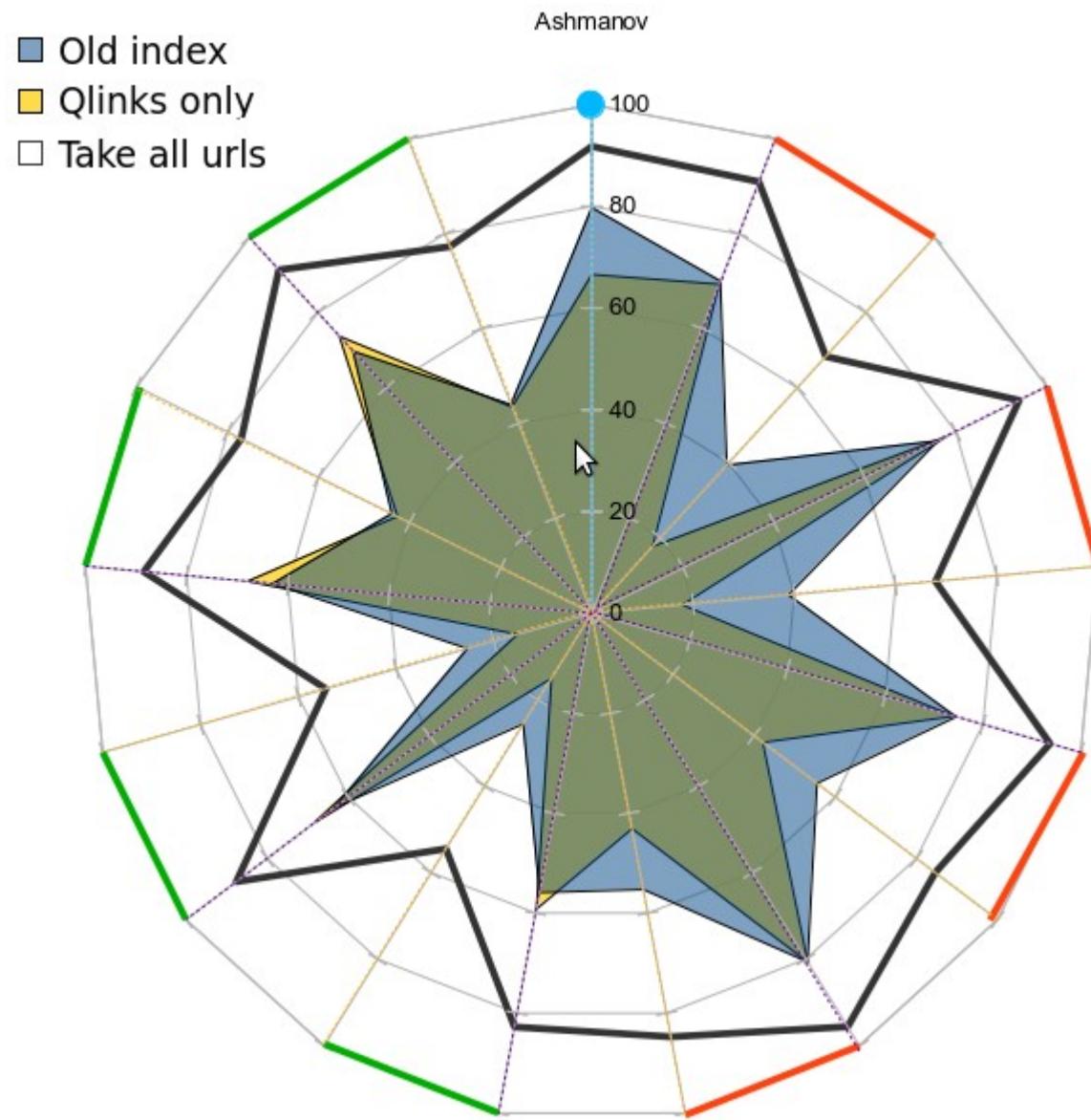






Index quality: baseline

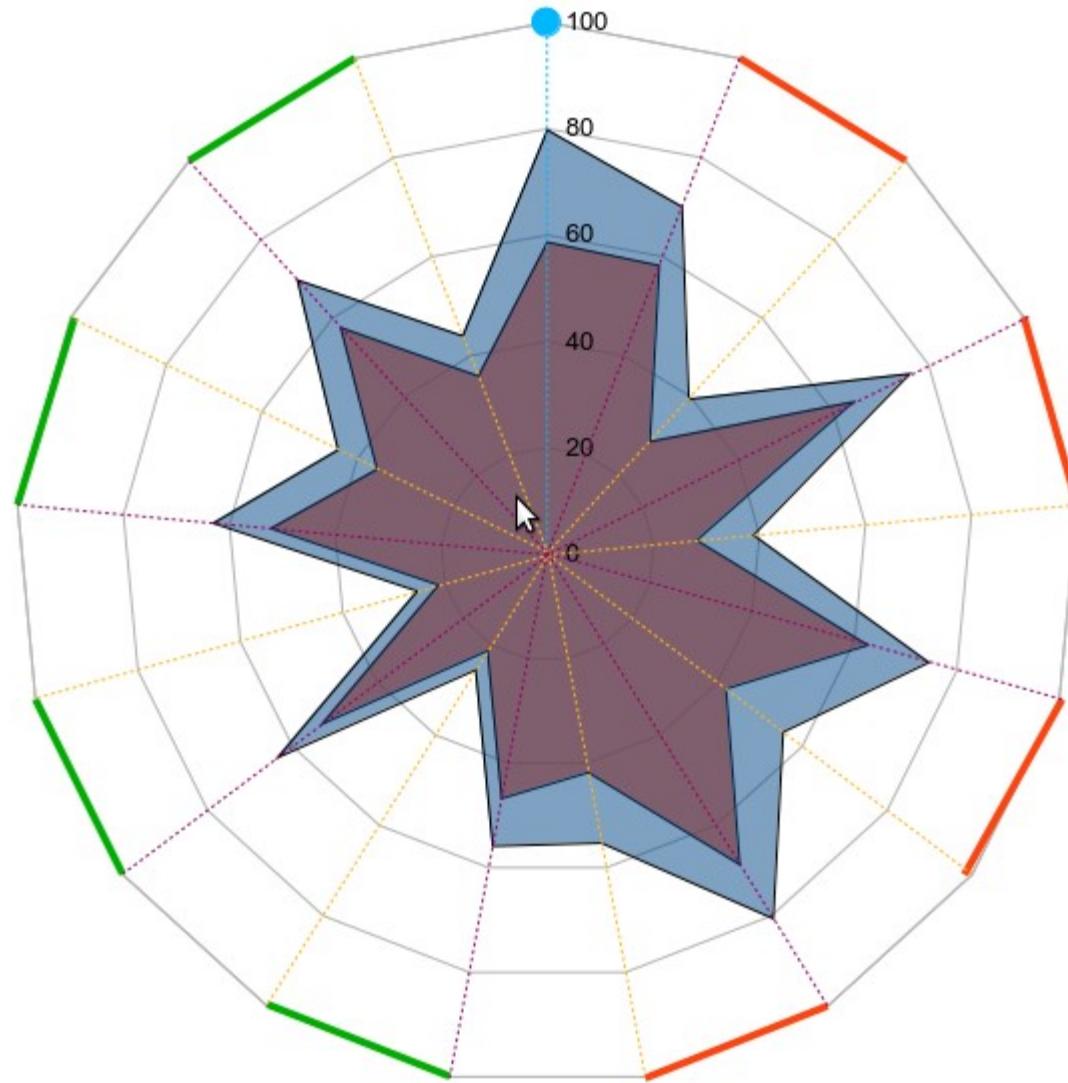
go.mail.ru



Index quality: old quota impact

■ Old index

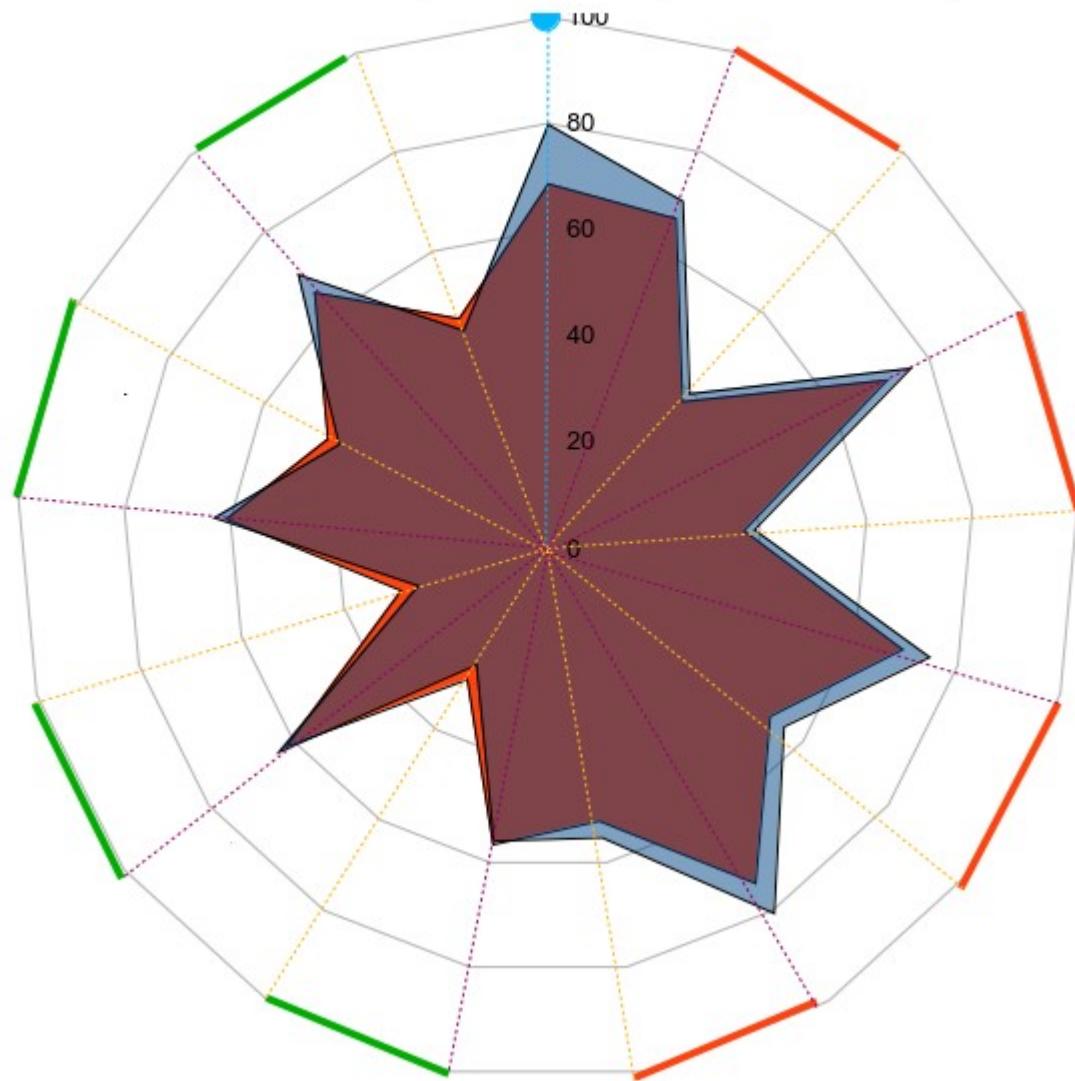
■ Old index, the same size, decreased quota



Index quality: new quota

■ Old index

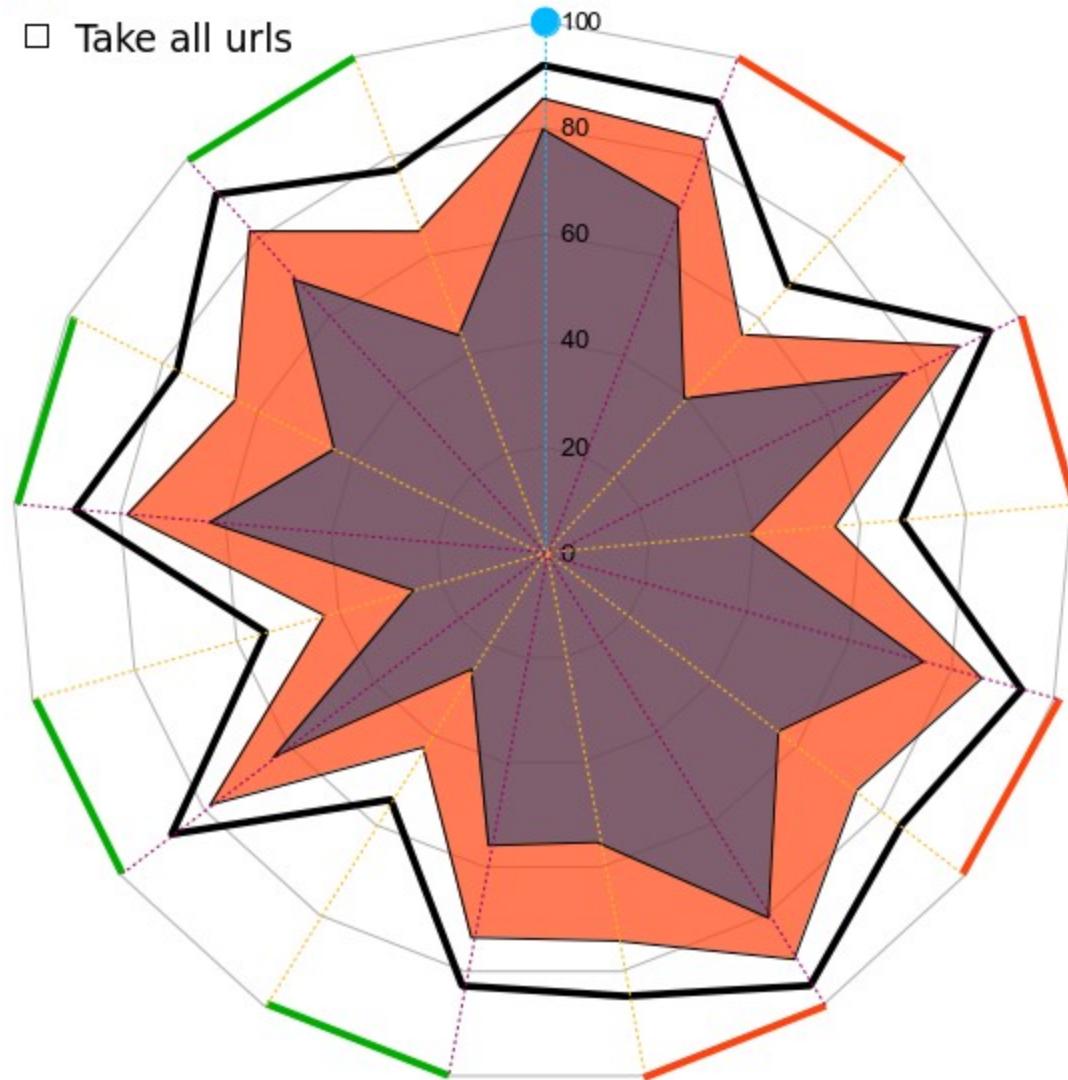
■ Old index cleaned by new quota (2x smaller size)



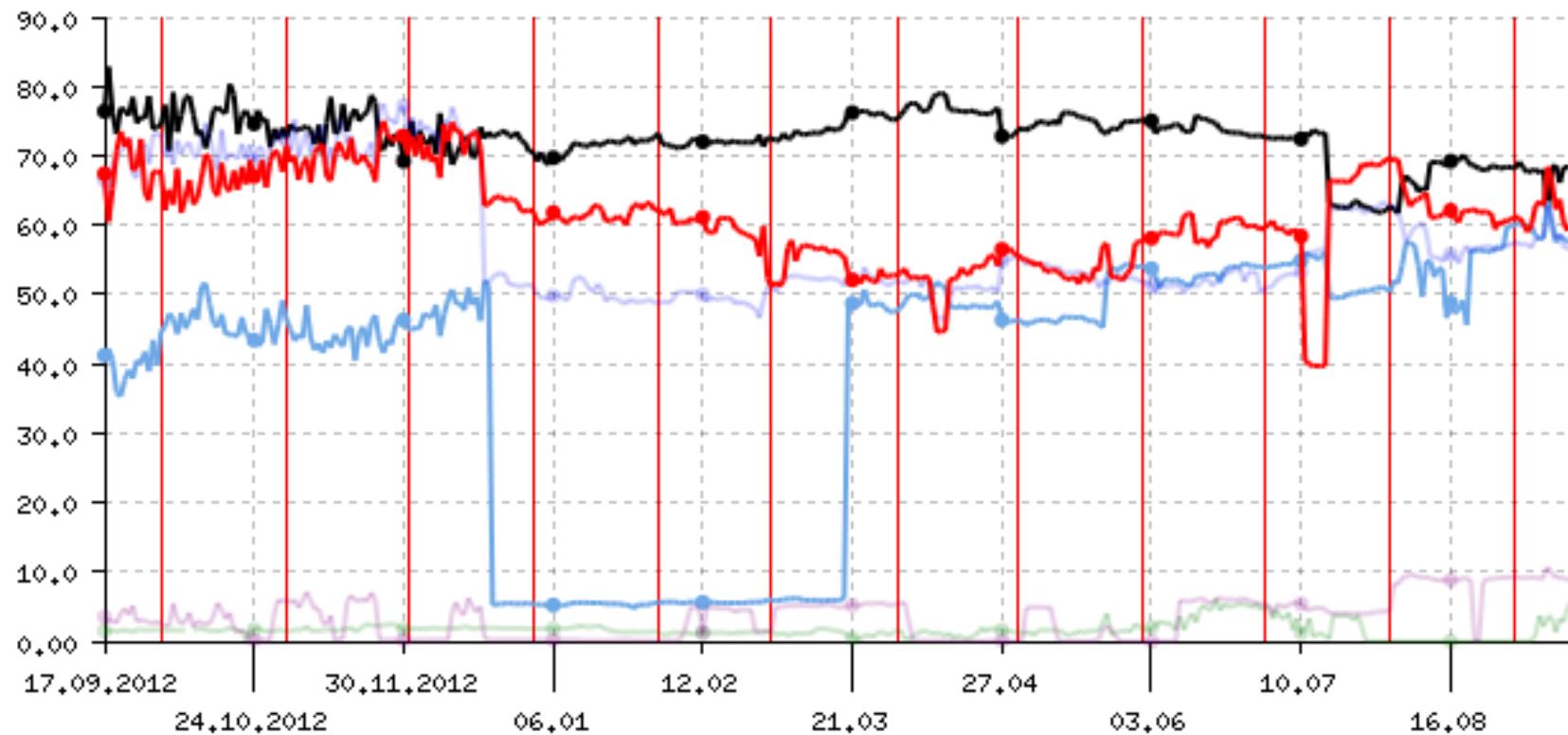
Index quality: final

go.mail.ru

- Old index
- New index
- Take all urls



recall



Mail
Yandex
Google

Static rank



№	Запрос	Bing (5.7 +0.8)	Google (68.6 +0.4)	@ Mail (57.4 -0.3)	Rambler (56.4 -0.4)	Yahoo (9.0 +0.1)	Яндекс (60.3 +1.1)
1	абузничество	0	1	4	1	0	1
2	авантюристовать	4 +4	16	17	17	4	17
3	антибомонд	6	40	28	19	6	19
4	аптайминг	0 -1	3	28	25	1	26
5	артпоповый	0	36 -2	15	19	1	20

Pin-Up Went Down - 342 (2010) - MetalArea

[metalarea.org](#) › ... › Releases Area / Релизы › Archive › Translate this page

Jul 9, 2010 - 30 posts - 16 authors

а альбом занятный. **артпоповый**, извините (с), такой. первый был концептуален, но я его переосмыслил и потер почти все, окромя:

Pin-Up Went Down - 342 (2010) - MetalArea

[metalarea.org](#) › ... › Releases Area / Релизы › Archive › Translate this page

Jul 9, 2010 - 30 posts - 16 authors

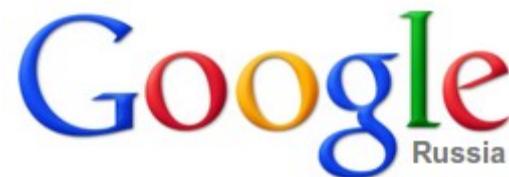
а альбом занятный. **артпоповый**, извините (с), такой. первый был концептуален, но я его переосмыслил и потер почти все, окромя:

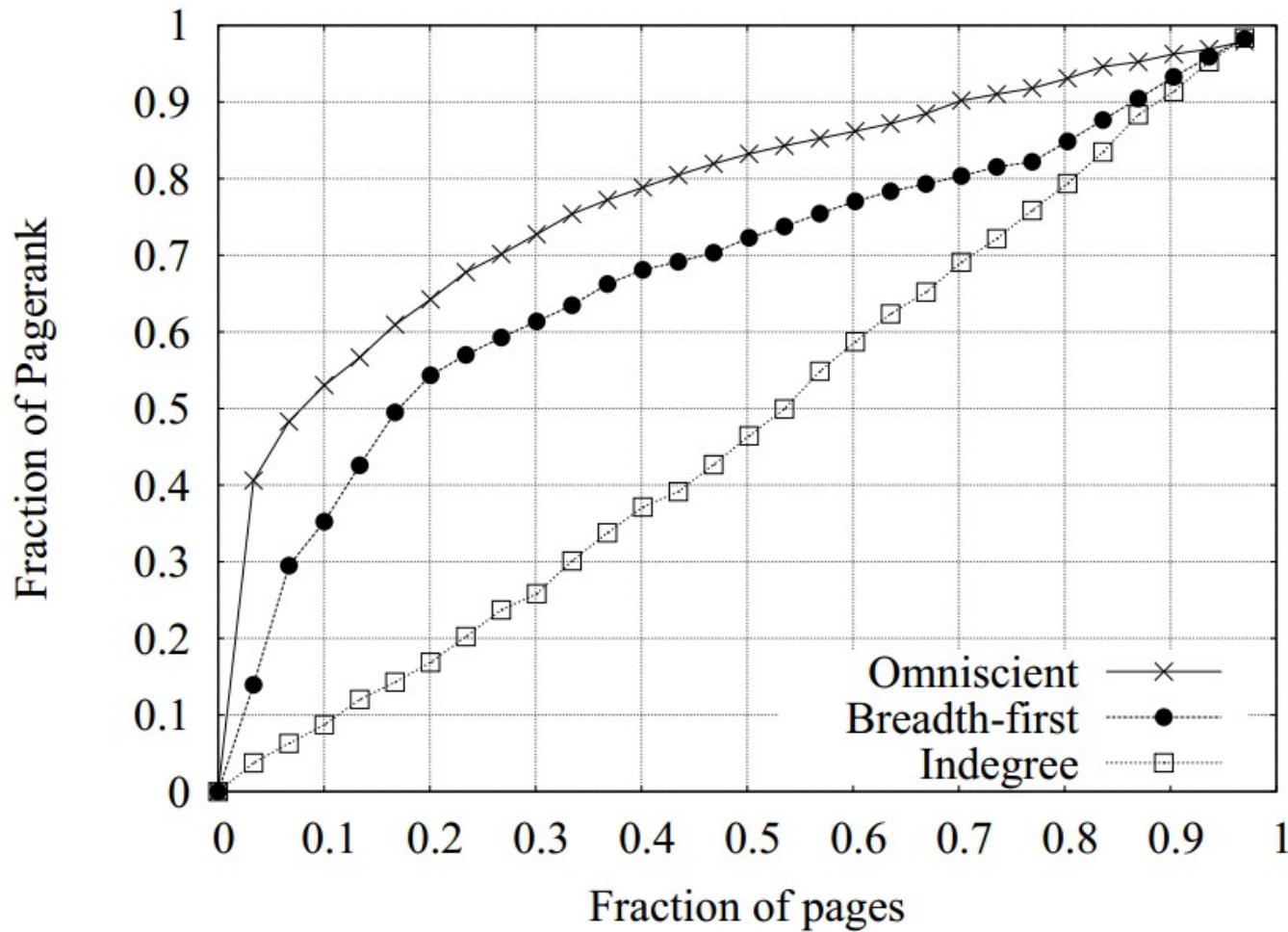
Pin-Up Went Down - 342 (2010) - MetalArea

[metalarea.org](#) › ... › Releases Area / Релизы › Archive › Translate this page

Jul 9, 2010 - 30 posts - 16 authors

а альбом занятный. **артпоповый**, извините (с), такой. первый был концептуален, но я его переосмыслил и потер почти все, окромя:



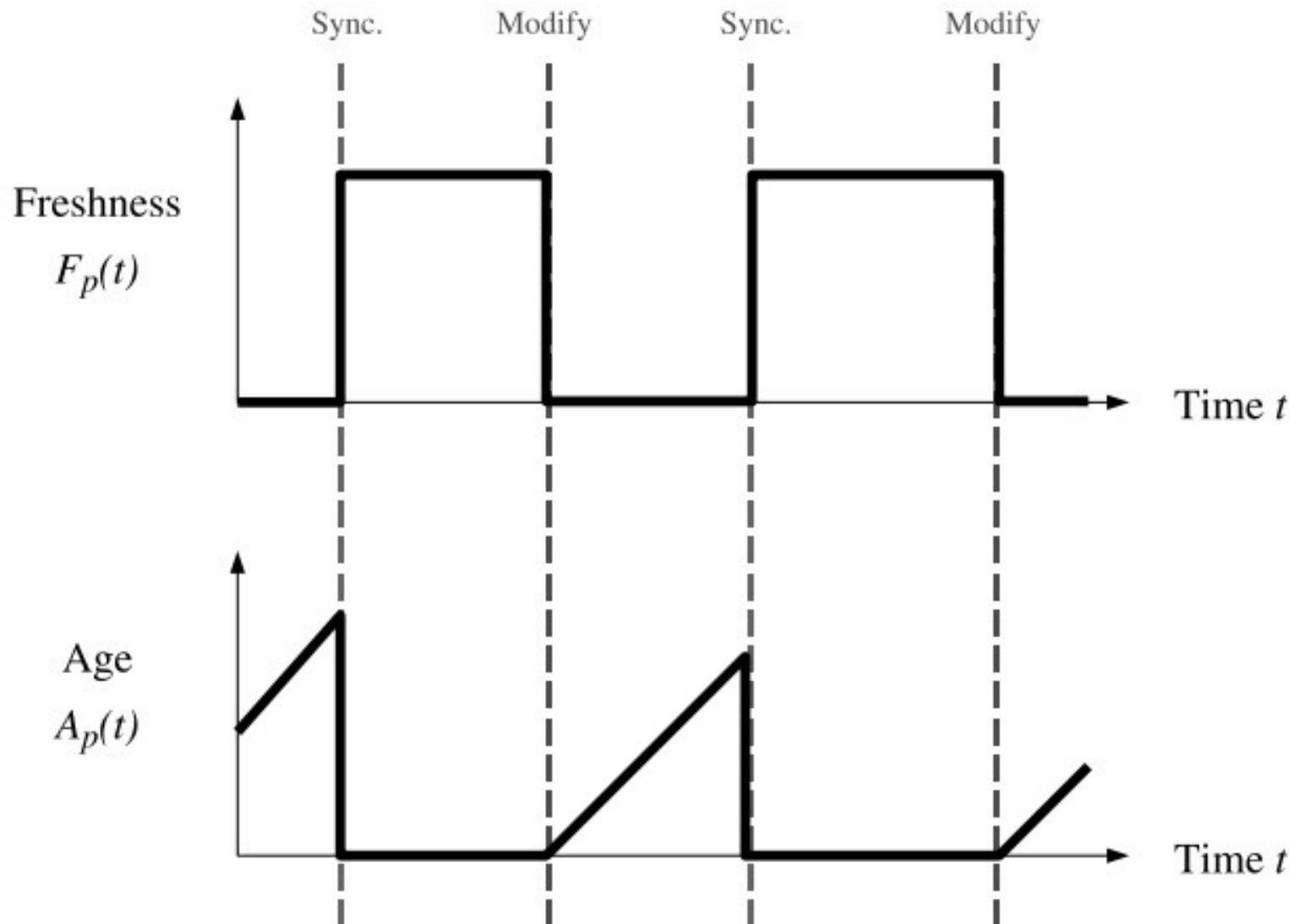


Events

- **Creation**
 - We use special discovery scheduler
- **Deletion**
 - Permanent vs temporal
 - Special issue with ban/robots
- **Update**
 - How to understand, that page was changed?
 - How to predict update?

Update prediction

go.mail.ru

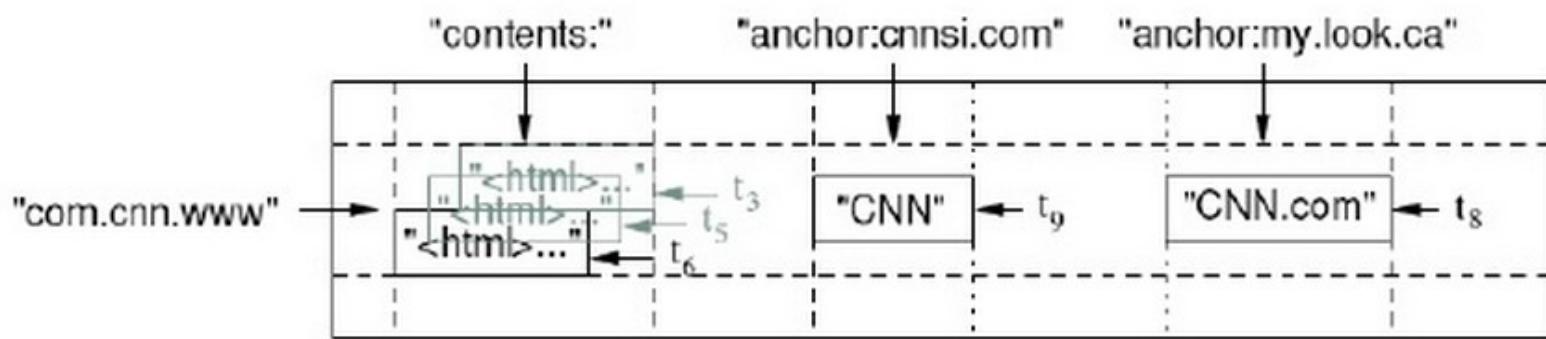




- Distributed file system (HDFS)
- MapReduce framework
- Hbase: no-sql database

**APACHE
HBASE**

- Doesn't support a full relational data model
- Multi-dimensional sorted map
- Indexed by a row, column, timestamp
 - (row: string, column: string, time : int 64) \rightarrow string



- Column-oriented storage
 - Most queries only involve a few columns out of many, so greatly reduces I/O.

Thank you!

Reference:

- **Ricardo Baeza-Yates.** Modern Information Retrieval:
The Concepts and Technology behind Search (2nd Edition), 2011