

Търсещи машини

Извличане на информация – Information Retrieval (IR)

- “Извличането на информация е област, занимаваща се със структурата, анализа, организацията, съхраняването, търсенето и извличането на информация” (Gerard Salton, 1968).
- Общата дефиниция може да бъде приложена към много типове информация и приложения за търсене.
- Основният фокус на IR от 1950 е върху текста и документите.

Какво е документ?

- Примери: web страници, е-мейли, книги, текстови съобщения, MS Word, PowerPoint, PDF, постинги, патенти и др.
- Общи свойства:
 - Значително текстово съдържание
 - Определена структура (заглавие, автор, тематика, подател и др.)

Документи и записи в бази данни

- Записите в базите данни (tuples) типично се състоят от добре дефинирани полета (attributes):
 - Банков запис с номера на сметки, баланси, имена, адреси, дата раждане и т.н.
- Лесни за сравняване полета с добре дефинирана семантика за заявките с цел откриване на съвпадение.
- Текстът е много по труден за обработка.

Документи и записи в бази данни - пример

- Пример за запитване към банкова база данни:
 - “Намери записи в баланс > 500000 във фирми, намиращи се в Лондон”;
 - Съвпаденията се намират лесно чрез сравнение със стойностите на полетата в записа.
- Пример за запитване към търсеща машина:
 - “Намери банковите скандали с България”;
 - Този текст трябва да бъде сравнен с текста във всички нови публикувани истории.

Сравняване на текст

- Сравняването на текста от запитването с текста в документа и определянето на добро съвпадение е ключов момент в IR.
- Не е необходимо точно съвпадение на думи:
 - Различни начини за запишем едно и също нещо с естествен език;
 - Някои публикации ще имат по-добро съвпадение от други.

Други медии

- Нови приложения, изискващи нови медии:
 - Видео, музика, изображения, говор.
- Съдържанието, както текста, е трудно за описание и сравнение:
 - Текст може да се използва за представяне на съдържание (тагове).

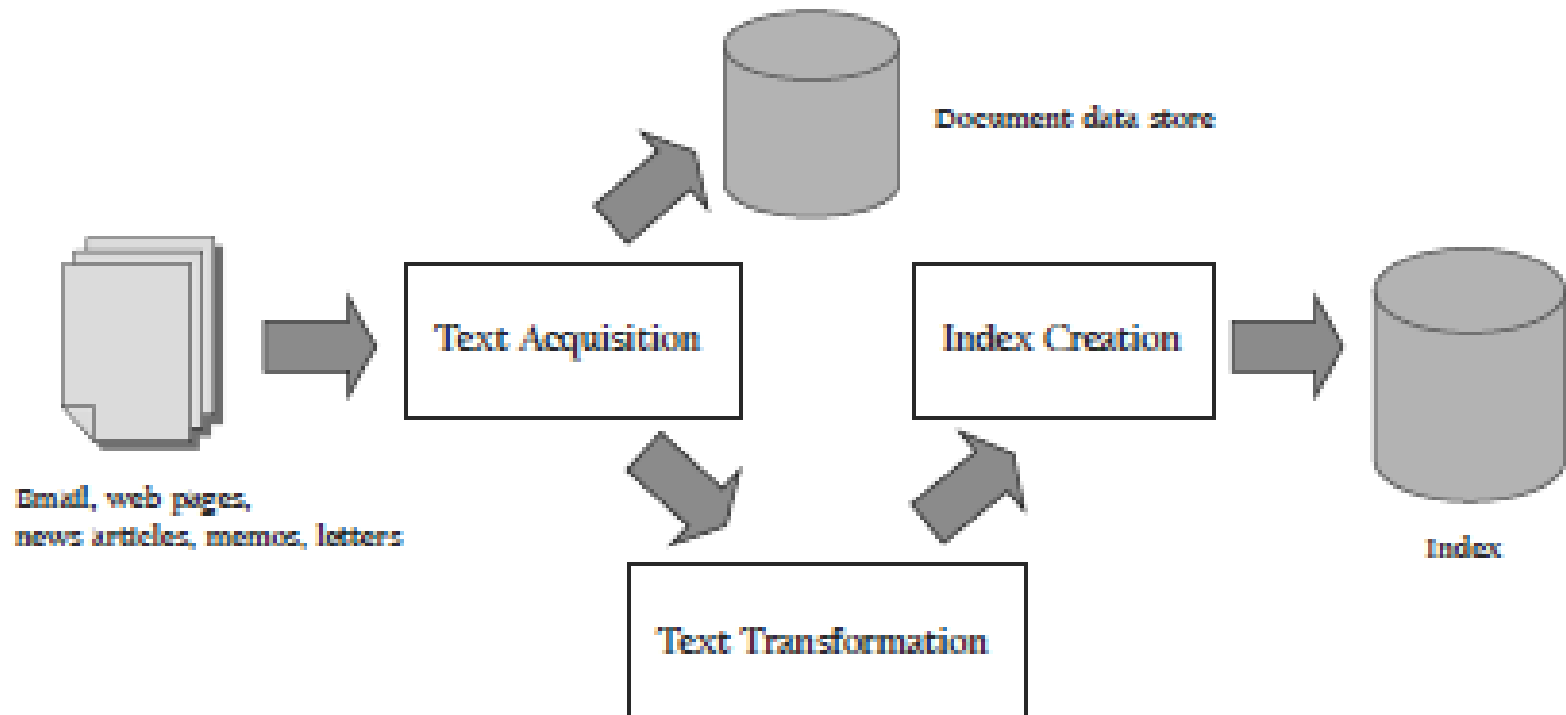
Архитектура на ТМ

- Софтуерната архитектура съдържа софтуерни компоненти, интерфейси, предоставени от компонентите и отношения между тях.
- Архитектурата се определя от 2 основни изисквания:
 - Ефикасност (време за отговор).
 - Ефективност (качество на резултатите).

Основни функции на ТМ

- Процес на индексиране.
- Процес на запитване.

Процес на индексиране



Изискване на текст

- Crawler:
 - Идентифицира и изисква документи за ТМ;
 - Различни типове – web, enterprise, desktop;
 - Web crawler-ите следват линковете за да откриват документи:
 - Трябва ефикасно да откриват огромен брой страници (coverage) и да ги поддържат актуални (freshness).

Изискване на текст

- Feeds
 - Потоци документи в реално време (блогове, жидео, радио, телевизия).
 - RSS reader е общ стандарт.
- Конвертиране
 - Конвертират се документи в консистентен текстов формат с метаданни (HTML, Word -> XML)
 - Конвертира се текстовото кодиране за различни езици.

Изискване на текст

- Хранилища на документи
 - Съхраняват текст, метаданни и друга съответстваща информация за документите
 - Метаданни – информация за документа (дата на създаване);
 - Друга информация (линкове, котви).
 - Предоставят бърз достъп до съдържанието на документите за всички компоненти на ТМ.
 - Могат да се използват релационни бази данни.

Трансформация на текст

- **Парсер**

- Обработване на последователност от текстови лексеми (tokens) за разпознаване на структурни елементи (заглавия, линкове и др.).
- Думите се разпознават от *Tokenizer*.
- Таговите езици (HTML, XML) често се използват за описание на структурата.

Трансформация на текст

- ***Stopping***
 - Премахване на общи думи.
 - Влияе на ефективността и ефикасността.
 - Може да е проблем при някои запитвания.

Трансформация на текст

- ***Stemming***
 - Групира думи, производни на общ корен.
 - Обикновено е ефективно, но не при всички запитвания.
 - Предимствата зависят от различните езици.

Трансформация на текст

- **Анализ на връзки (*Link analysis*)**
 - Използва линковете и котвите в WEB страниците.
 - Анализът определя *популярността* (PageRank).
 - Текстът в котвите може значително да разшири обхвата на страниците, сочени от линковете.
 - Има значително влияние при търсенето в WEB.

Трансформация на текст

- ***Извличане на информацията***
 - Идентифицира класове на индексни термини, които са важни за определени приложения (класове като хора, компании, дати)

Трансформация на текст

- **Класификатор**
 - Идентифицира класово-базирани метаданни за документите (присвоява етикети към документи)
 - Зависи от конкретното приложение.

Създаване на индекс

- **Статистика за документите**
 - Събира се брой и позиция на думите и други свойства
 - Използва се в алгоритмите за рейтинговане.

Създаване на индекс

- ***Задаване на тегла (Weighting)***
 - Изчислява тегла за индексните термини.
 - Използва се в алгоритмите за рейтинговане

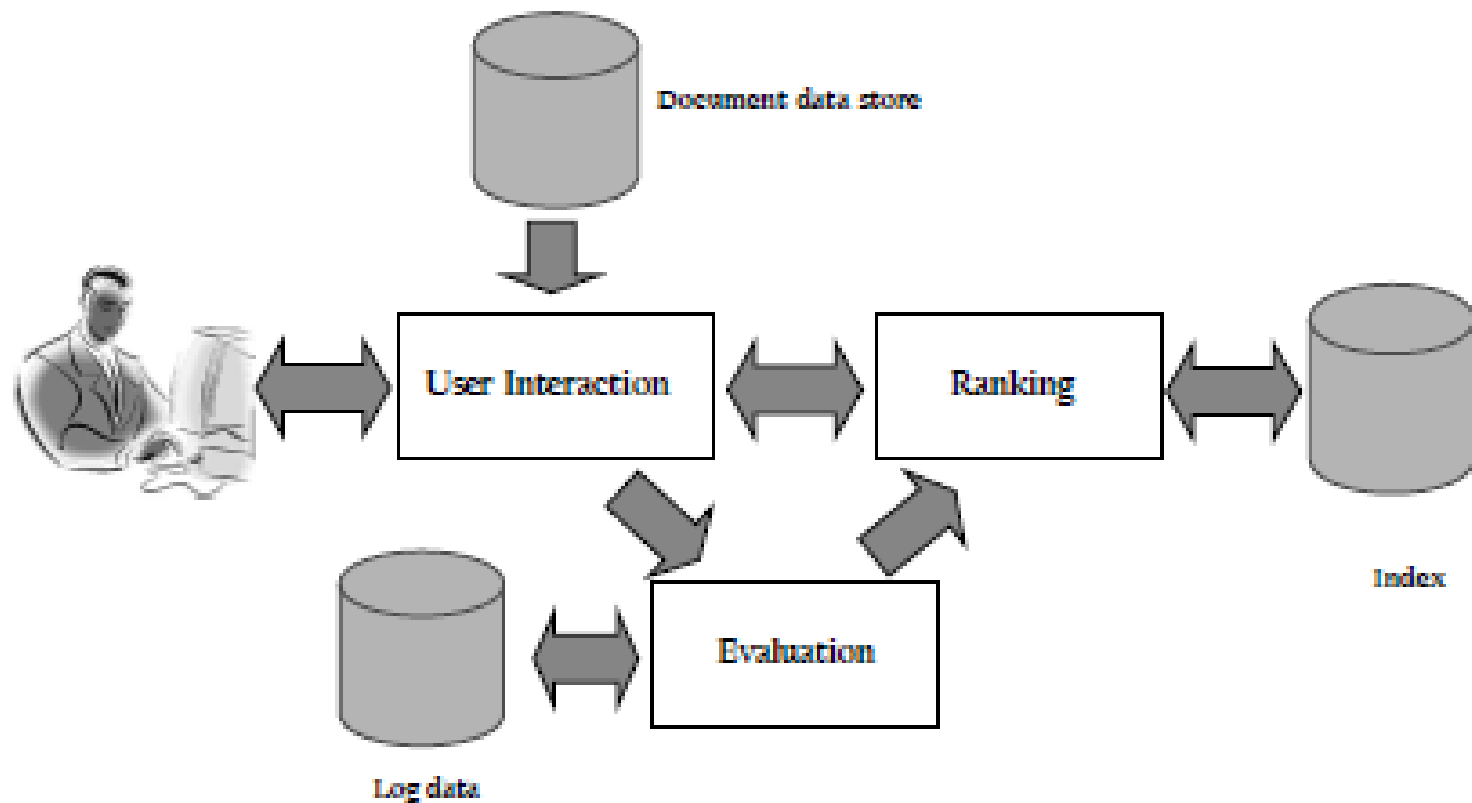
Създаване на индекс

- **Инвертиране**
 - Същината на процеса на индексирание.
 - Преобразува информацията от вида документи-термини във вида термини-документи.
 - Форматът на инвертния файл е проектиран за бърза обработка на запитвания.

Създаване на индекс

- ***Разпределяне на индекса***
 - Разпределяне на индексите между множество компютри и/или сайтове.
 - Изключително важност за бърза обработка на запитвания с огромен брой документи.
 - Вариации (разпределян на документи, на термини, репликации).

Процес на запитване



Взаимодействие с потребителя

- ***Въвеждане на запитвания***
 - Предоставя интерфейс и парсер на езика за запитвания.
 - Повечето web запитвания са много опростени (кавички).
 - Някои приложения могат да използват форми.
 - За описание на сложни запитвания се използват езици за запитвания (query languages).

Взаимодействие с потребителя

- **Трансформация на запитванията**
 - Подобрява началното запитване (и преди и след първоначалното търсене).
 - Използват се техники за трансформиране на текст, както при документите.
 - *Spell checking, query suggestion* предоставят алтернатива на оригиналните запитвания.
 - *Query expansion, relevance feedback* модифицират оригиналното запитване с допълнителни термини.

Взаимодействие с потребителя

- **Извеждане на резултати**
 - Формиране на извеждане на рейтинговани документи за запитването.
 - Генериране на *snippets* за показване как запитванията съответстват на документите.
 - Подчертаване на важни думи и пасажии.
 - Получаване на подходящи реклами в много приложения.

Рейтинговане (Ranking)

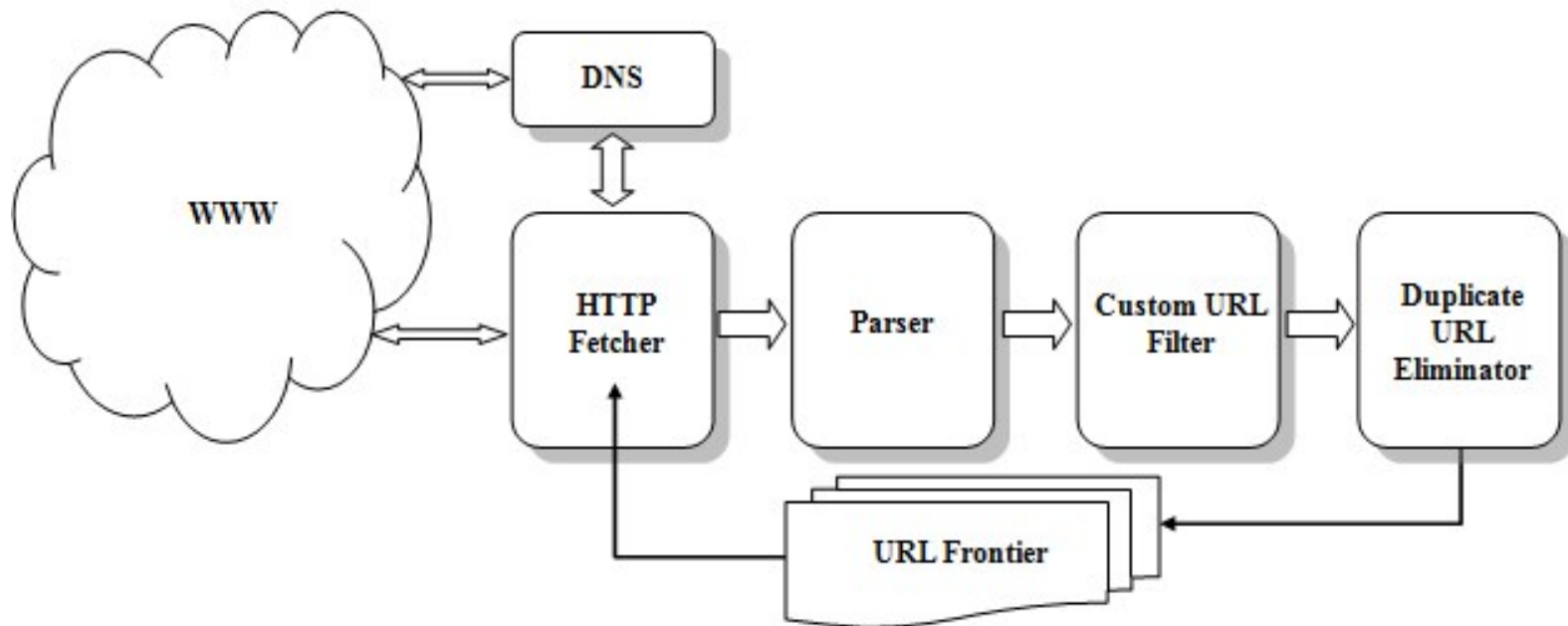
- **Формиране на оценка**
 - Изчисляват се оценки за документите на базата на алгоритми за рейтинговане.
 - Ключов компонент на всяка ТМ.
 - Основен вид на оценката е

$\sum q_i d_i$ — теглата на термините на запитването и на документа за термин i .

Оценяване

- ***Logging***
 - Съхраняване на потребителските запитвания и взаимодействие (*clickthrough data, dwell time*).
- ***Анализ на рейтинговането***
 - Измерване и настройване на ефективността на рейтинговането.
- ***Анализ на производителността***
 - Измерване и настройване на ефективността на система за търсене (тестови колекции).

Архитектура на Web crawler



Особености

- Налице е изчакване за получаване на отговора.
- Редуцирането на това време – нишки.
- Потенциално засипват сайтовете със заявки за страници (flooding).
- Избягване на това – изчакване между заявките към един и същ web сървър (*politeness policies*).

Пример за нишка на crawler

```
procedure CRAWLERTHREAD(frontier)
  while not frontier.done() do
    website ← frontier.nextSite()
    url ← website.nextURL()
    if website.permitsCrawl(url) then
      text ← retrieveURL(url)
      storeDocument(url, text)
      for each url in parse(text) do
        frontier.addURL(url)
      end for
    end if
    frontier.releaseSite(website)
  end while
end procedure
```


Актуалност на страниците (freshness)

- Използва се специална заявка на протокола HTTP – HEAD.

Client request: HEAD /csinfo/people.html HTTP/1.1
Host: www.cs.umass.edu

HTTP/1.1 200 OK
Date: Thu, 03 Apr 2008 05:17:54 GMT
Server: Apache/2.0.52 (CentOS)
Last-Modified: Fri, 04 Jan 2008 15:28:39 GMT
Server response: ETag: "239c33-2576-2a2837c0"
Accept-Ranges: bytes
Content-Length: 9590
Connection: close
Content-Type: text/html; charset=ISO-8859-1

Deep Web

- Трудни за извличане web сайтове:
 - Частни сайтове;
 - Резултати от форми;
 - Скриптови страници.

Обработка на текст

- Парсване на текста (*parsing*) – разпознаване на съдържанието и структурата на документите.
- Tokenizing (lexical analysis) – формиране на думи.
- Съдържание на документ:
 - Думи (tokens);
 - Метаданни (автор, тагове).
- Синтактичен анализ

Проблеми

- Кратки думи могат да бъдат важни в някои запитвания (*xp*, *pt*, *world war II*).
- Има думи с и без тирета:
 - В някои случаи не са необходими (*e-bay*, *active-x*, *cd-rom*).
 - В други са част от думата или разделител (*mazda rx-7*, *e-cards*, *t-mobile*)

Проблеми

- Специални символи, важни за тагове, URL и код.
- Главни букви, имащи различно значение (*Bush*, *Apple*).
- Апострофите могат да са част от дума, притежание или грешка (*can't*, *80's*, *o'donnell*).

Проблеми

- Числата имат важно значение (*Porsche 911, top 10 courses, Windows 10*).
- Точките могат да се появяват в числа, аббревиатури, URL, край на изречение (*I.B.M., Ph.D. , cs.tu-varna.bg*)

Stemming

- Алгоритмични
- Базирани на речници

- Suffix-s stemmer

cakes -> cake

dogs -> dog

Но: *century – centuries ?*

Porter stemmer

- Последователност от стъпки
- На всяка стъпка се премахват или заменят окончания.
- Заменя sses с ss (*stresses* -> *stress*).

Krovetz stemmer

- Хибриден подход.
- Използва речник за определяне дали думата е валидна.
- Помощни списъци с деривационни суфикси.
- Предимство – корените в повечето случай са пълни думи.

Stemmers

Original text:

Document will describe marketing strategies carried out by U.S. companies for their agricultural chemicals, report predictions for market share of such chemicals, or report market statistics for agrochemicals, pesticide, herbicide, fungicide, insecticide, fertilizer, predicted sales, market share, stimulate demand, price cut, volume of sales.

Porter stemmer:

document describ market strategi carri compani agricultur chemic report predict market share chemic
report market statist agrochem pesticid herbicid fungicid insecticid fertil predict sale market share
stimul demand price cut volum sale

Krovetz stemmer:

document describe marketing strategy carry company agriculture chemical report prediction market
share chemical report market statistic agrochemic pesticide herbicide fungicide insecticide fertilizer
predict sale stimulate demand price cut volume sale

Други езици

kitab	<i>a book</i>
kitabı	<i>my book</i>
alkıtab	<i>the book</i>
kitabıki	<i>your book (f)</i>
kitabıka	<i>your book (m)</i>
kitabıhu	<i>his book</i>
kataba	<i>to write</i>
maktaba	<i>library, bookstore</i>
maktab	<i>office</i>

ktb

Изграждане на индекси

- Проектирани са да поддържат *търсене*.
- Търсещите машини използва специфична форма на търсене – *ranking*
 - Документите се извличат в подреден ред в зависимост от оценката на документа.

Текст за индексиране

d1: The jaguar is a New World mammal of the Felidae family.

d2: Jaguar has designed four new engines.

d3: For Jaguar, Atari was keen to use a 68K family device.

d4: The Jacksonville Jaguars are a professional US football team.

d5: Mac OS X Jaguar is available at a price of US \$199 for Apple's new "family pack".

d6: One such ruling family to incorporate the jaguar into their name is Jaguar Paw.

d7: It is a big cat.

Отделяне на думите

d1: the₁ jaguar₂ is₃ a₄ new₅ world₆ mammal₇ of₈ the₉ felidae₁₀ family₁₁

d2: jaguar₁ has₂ designed₃ four₄ new₅ engines₆

d3: for₁ jaguar₂ atari₃ was₄ keen₅ to₆ use₇ a₈ 68k₉ family₁₀ device₁₁

d4: the₁ jacksonville₂ jaguars₃ are₄ a₅ professional₆ us₇ football₈ team₉

d5: mac₁ os₂ x₃ jaguar₄ is₅ available₆ at₇ a₈ price₉ of₁₀ us₁₁ \$199₁₂ for₁₃ apple's₁₄ new₁₅ family₁₆ pack₁₇

d6: one₁ such₂ ruling₃ family₄ to₅ incorporate₆ the₇ jaguar₈ into₉ their₁₀ name₁₁ is₁₂ jaguar₁₃ paw₁₄

d7: it₁ is₂ a₃ big₄ cat₅

Нормализиране на думи

d1: the₁ jaguar₂ be₃ a₄ new₅ world₆ mammal₇ of₈ the₉ felidae₁₀ family₁₁

d2: jaguar₁ have₂ design₃ four₄ new₅ engine₆

d3: for₁ jaguar₂ atari₃ be₄ keen₅ to₆ use₇ a₈ 68k₉ family₁₀ device₁₁

d4: the₁ jacksonville₂ jaguars₃ be₄ a₅ professional₆ us₇ football₈ team₉

d5: mac₁ os₂ x₃ jaguar₄ be₅ available₆ at₇ a₈ price₉ of₁₀ us₁₁ \$199₁₂ for₁₃ apple₁₄ new₁₅ family₁₆ pack₁₇

d6: one₁ such₂ rule₃ family₄ to₅ incorporate₆ the₇ jaguar₈ into₉ their₁₀ name₁₁ be₁₂ jaguar₁₃ paw₁₄

d7: it₁ be₂ a₃ big₄ cat₅

Премахване на общите термини

d1: jaguar₂ new₅ world₆ mammal₇ felidae₁₀ family₁₁

d2: jaguar₁ design₃ four₄ new₅ engine₆

d3: jaguar₂ atari₃ keen₅ 68k₉ family₁₀ device₁₁

d4: jacksonville₂ jaguars₃ professional₆ us₇ football₈ team₉

*d5: mac₁ os₂ x₃ jaguar₄ available₆ price₉ us₁₁ \$199₁₂ apple₁₄ new₁₅
family₁₆ pack₁₇*

*d6: one₁ such₂ rule₃ family₄ incorporate₆ jaguar₈ their₁₀ name₁₁ jaguar₁₃
paw₁₄*

d7: big₄ cat₅

Индекси

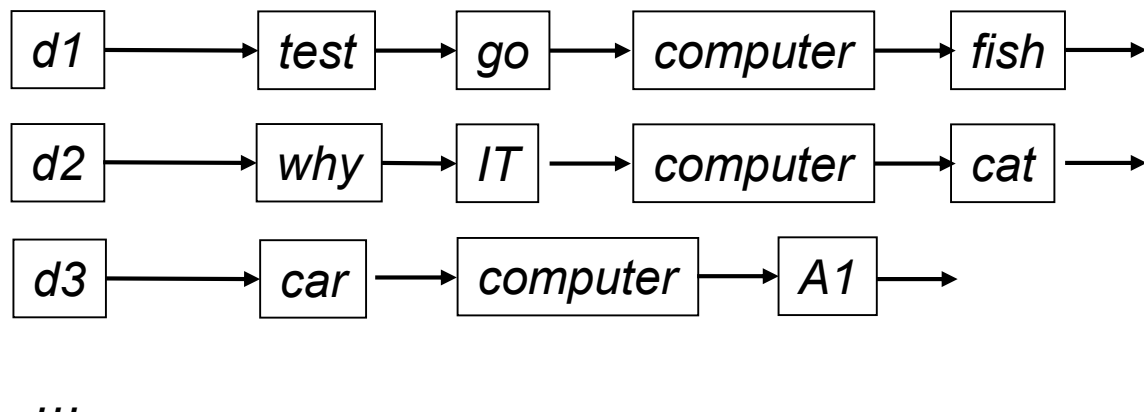
Forward index – списък от документи, във всеки документ се съдържат думи.

документи -> към -> думи

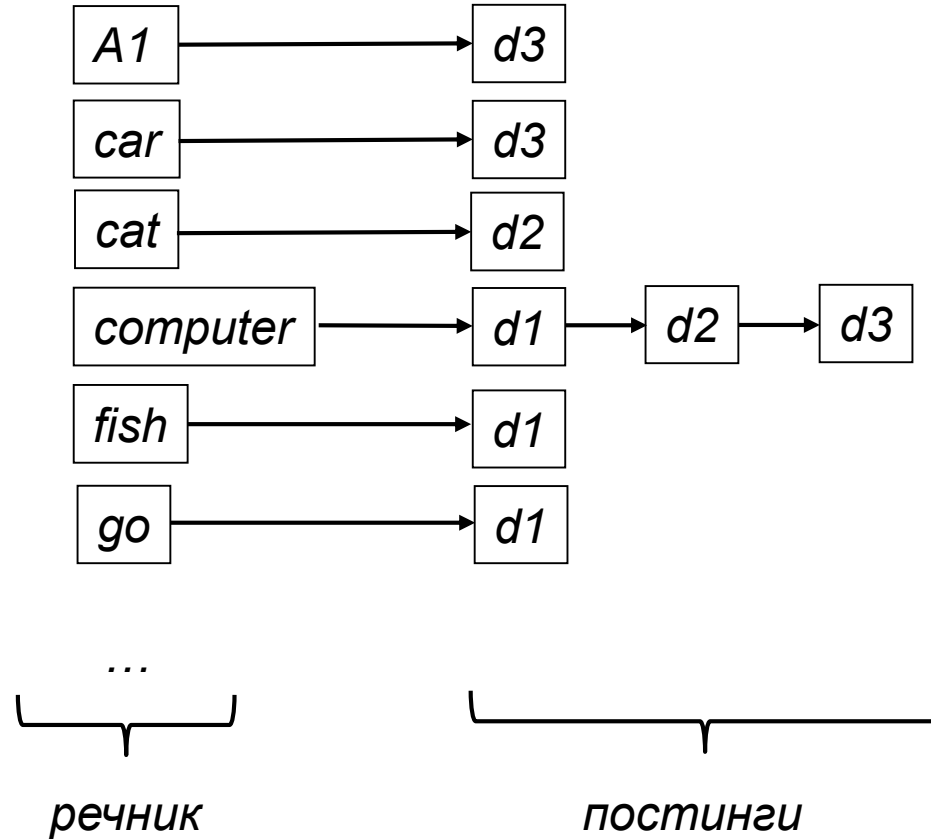
Inverted index – списък от думи и документи в които те се съдържат.

думи -> към -> документи

Forward index



Inverted index



Инвертен индекс - пример

d1: jaguar₂ new₅ world₆ mammal₇ felidae₁₀ family₁₁

d2: jaguar₁ design₃ four₄ new₅ engine₆

d3: jaguar₂ atari₃ keen₅ 68k₉ family₁₀ device₁₁

d4: jacksonville₂ jaguars₃ professional₆ us₇ football₈ team₉

d5: mac₁ os₂ x₃ jaguar₄ available₆ price₉ us₁₁ \$199₁₂ apple₁₄
new₁₅ family₁₆ pack₁₇

d6: one₁ such₂ rule₃ family₄ incorporate₆ jaguar₈ their₁₀
name₁₁ jaguar₁₃ paw₁₄

d7: big₄ cat₅

family
football
jaguar
new
rule
us
world
...

d_1, d_3, d_5, d_6

d_4

$d_1, d_2, d_3, d_4, d_5, d_6$

d_1, d_2, d_5

d_6

d_4, d_5

d_1

Позиция на думата

<i>family</i>	$d_1/11, d_3/10, d_5/16, d_6/4$
<i>football</i>	$d_4/8$
<i>jaguar</i>	$d_1/2, d_2/1, d_3/2, d_4/3, d_5/4, d_6/8$
<i>new</i>	$d_1/5, d_2/5, d_5/15$
<i>rule</i>	$d_6/3$
<i>us</i>	$d_4/7, d_5/11$
<i>world</i>	$d_1/6$
<i>...</i>	

Релевантност

- **Дефиниция:** *Релевантния документ съдържа информацията, която потребителят търси, когато изпраща заявка към търсещата машина.*
- Редица фактори оказват влияние на преценката на потребителя кое е релевантно – съдържание, новост, стил и др.
- **Тематична релевантност** (същата тема) и **Потребителска релевантност** (всичко допълнително).

Оценка релевантността на документите

- Релевантността се измерва чрез присвояване на тегло към срещанията на термина в документа.
- Често срещан термин в документа е повече релевантен за индексирание на документа.
- Рядко срещан термин в колекция засилва релевантността на документ.
- Често срещан термин в много документи е по-малко отличителен.

Оценка релевантността на документите

term frequency - inverse document frequency (tf-idf)

- **term frequency** – честотата на срещанията на термин към общия брой термини в документ

$$tf(t, d) = \frac{n_{t,d}}{\sum_{t'} n_{t',d}}$$

Оценка релевантността на документите

term frequency - inverse document frequency (tf-idf)

- **inverse document frequency** – важността на термин в колекция от документи

$$idf(t) = \log \frac{|D|}{|\{d' \in D | n_{t,d'} > 0\}|}$$

Оценка релевантността на документите

term frequency - inverse document frequency (tf-idf)

tf - idf

$$tfidf(t, d) = \frac{n_{t,d}}{\sum_{t'} n_{t',d}} \cdot \log \frac{|D|}{|\{d' \in D | n_{t,d'} > 0\}|}$$

Модифициран индекс

family	$d_1/11/.13, d_3/10/.13, d_5/16/.7, d_6/4/.8$
football	$d_4/8/.47$
jaguar	$d_1/2/.04, d_2/1/.04, d_3/2/.04, d_4/3/.04, d_5/4/.04, d_6/8/.04$
new	$d_1/5/.20, d_2/5/.24, d_5/15/.10$
rule	$d_6/3/.28$
us	$d_4/7/.30, d_5/11/.15$
world	$d_1/6/.47$
...	

Размер на индекса

Пример:

Колекция от е-мейли. Всеки мейл е с размер средно 1Kb и всеки мейл съдържа средно по 100 думи.

Колекцията съдържа 1 000 000 мейла в обем 1Gb.

Брой думи: 100×10^6

След парсване и токенизиране: 200 000 различни термина

1. Индексът съдържа 200 000 списъка
2. При допускане, че 20% от термините в документ се срещат 2 пъти, то всеки документ се среща в 80 списъка
3. Всеки списък съдържа средно 400 записа
4. При представяне на ID на документ с 4byte integer, то средния размер на списък е 1600bytes
5. Целият индекс ще съдържа $400 \times 200\,000 = 80\,000\,000$ записа
6. **Размерът на индекса е 320 Mb**

Размер на индекса

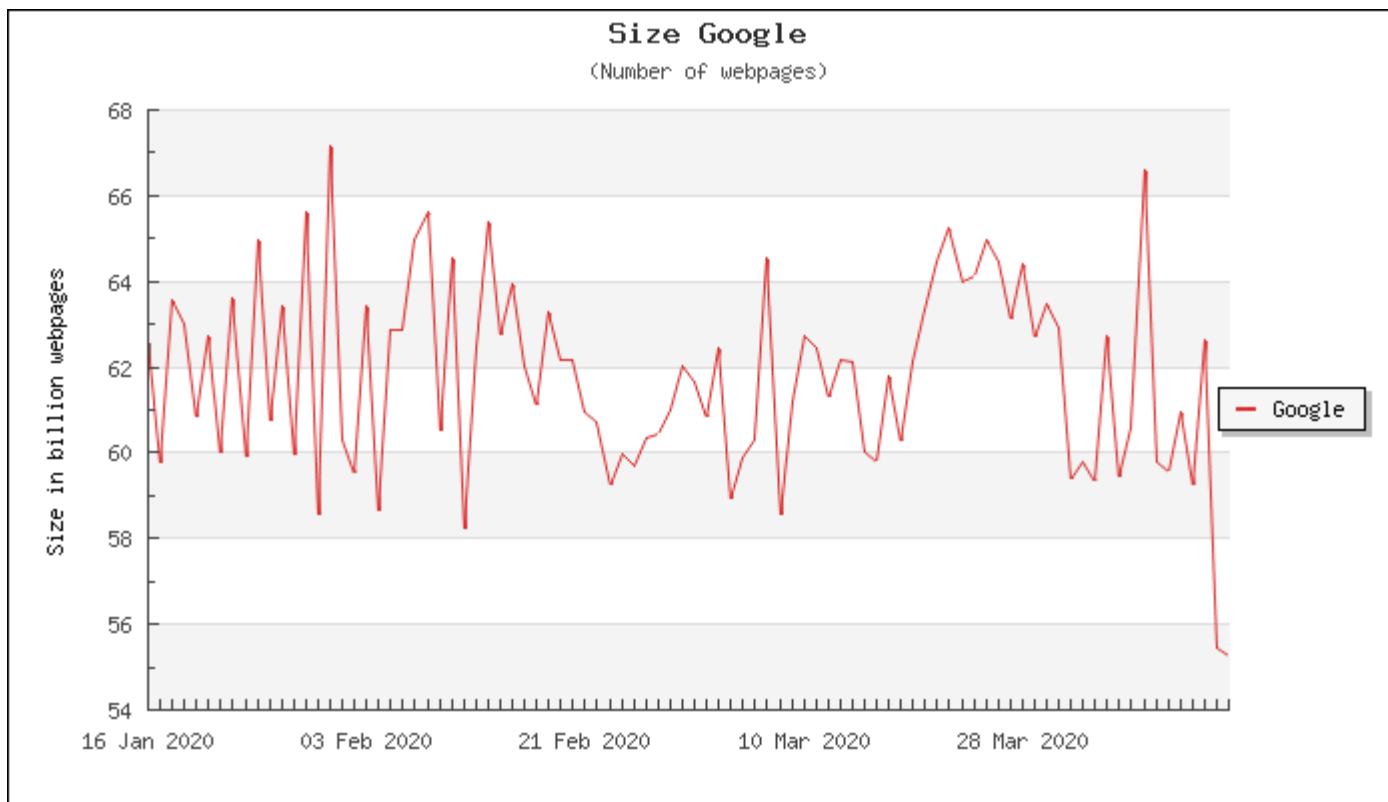
Добавяне на допълнителна информация към елементите в списъците:

- Позиция
- Тегло

Допълнително заета памет: 8 bytes

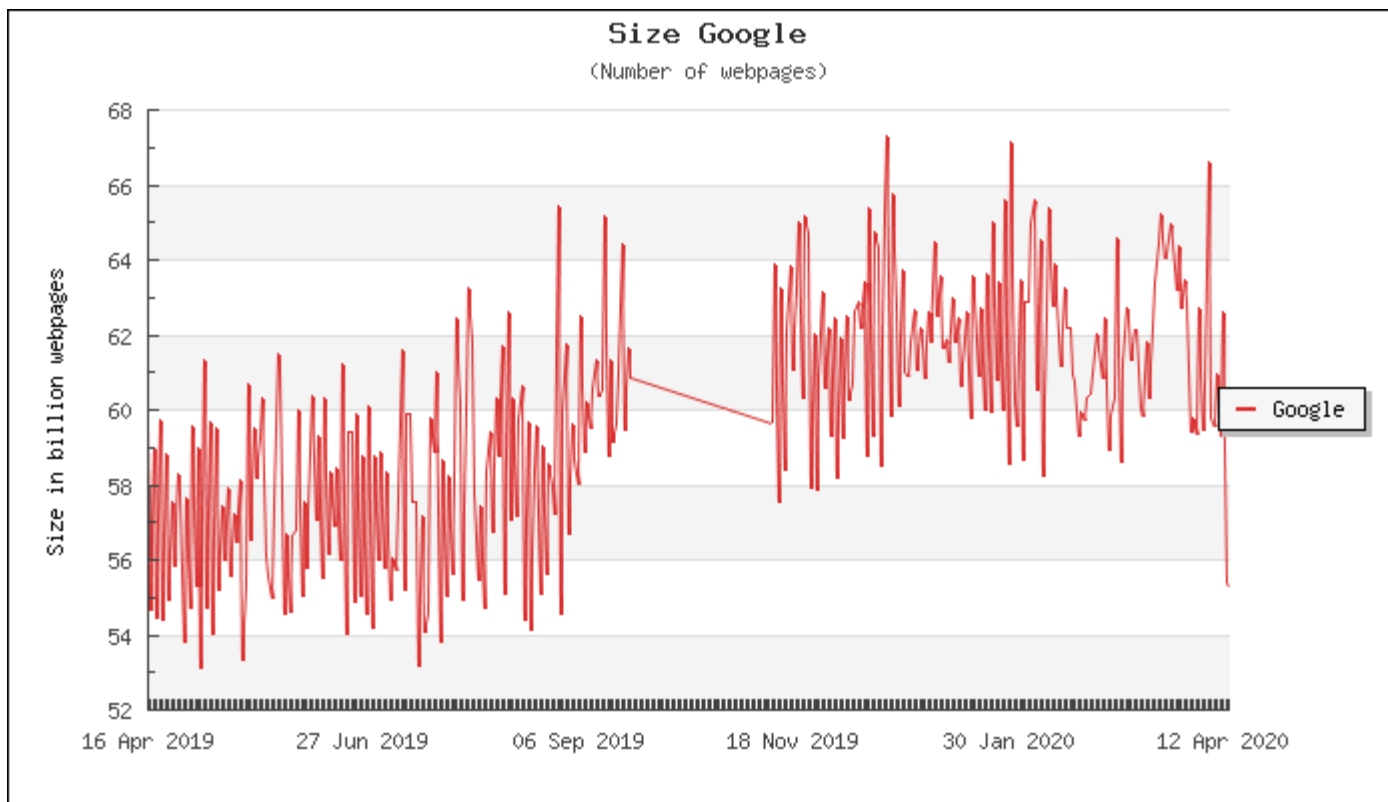
Целият индекс ще съдържа $80 \times 12 \times 10^6 = 960 \text{ Mb}$

Индексиране в WWW

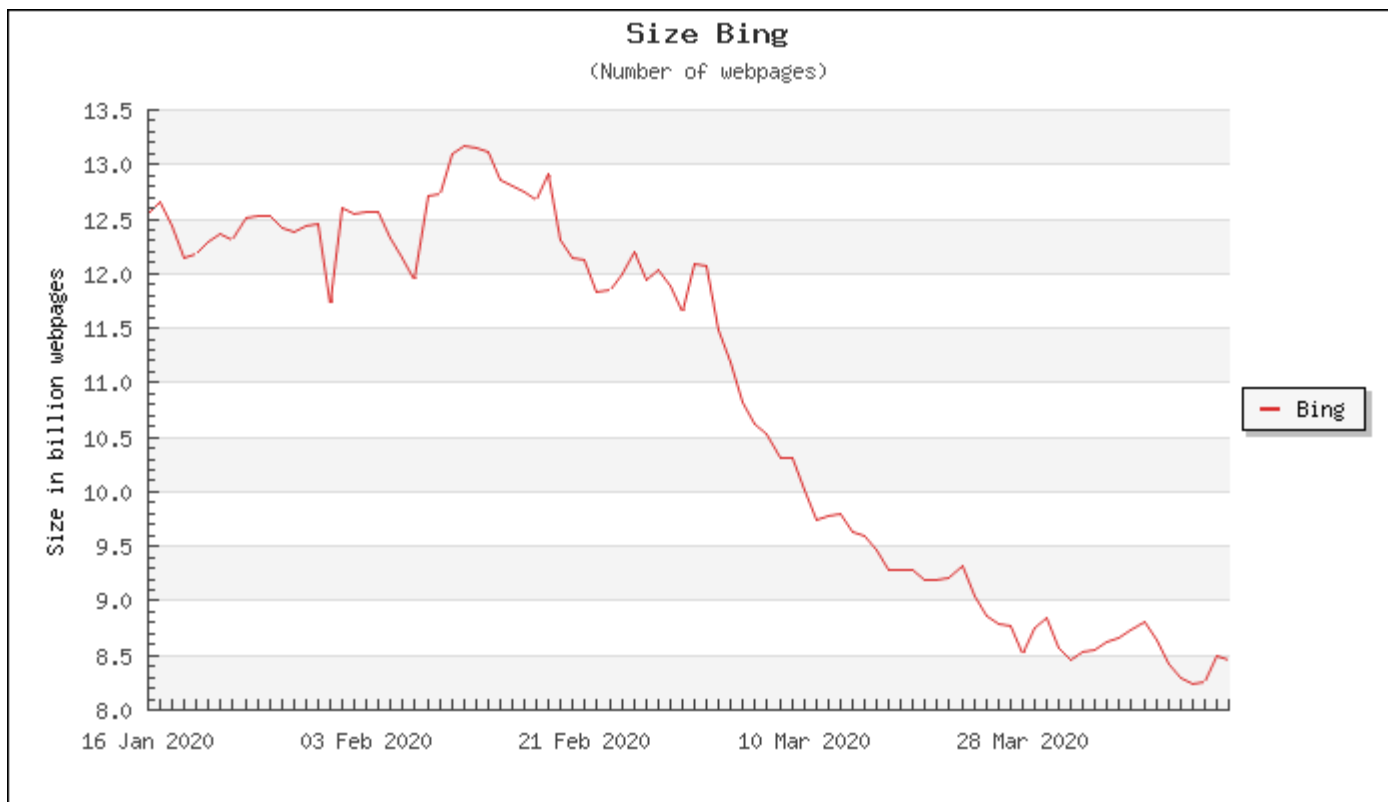


<https://www.worldwidewebsize.com/>

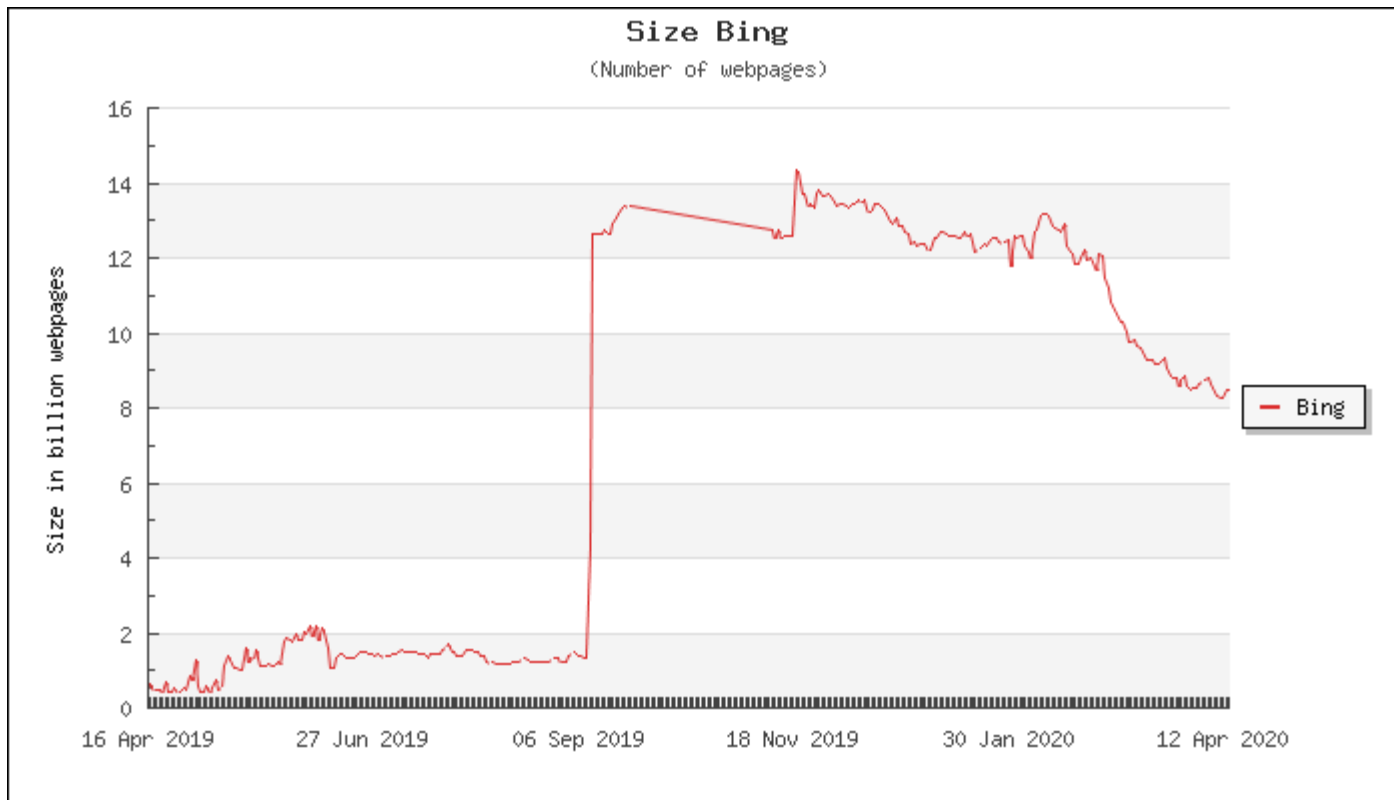
Индексиране в WWW



Индексиране в WWW



Индексиране в WWW



Анализ на връзките

- Линковете са ключов компонент на Web.
- Важни са за навигацията, но и за търсенето.

Линкове

`Example website`

Оценка на връзките

- Милиарди страници, някои повече информативни от други;
- Линковете могат да се разглеждат като информация за популярността на web страница;
- Броят идващи линкове (*inlinks*) може да се използва като проста оценка;
- Алгоритми за анализ на връзките.

Page Rank


- Базира се на броя на „идващите“ линкове към страницата и популярността на страницата (weight) от която те идват.
- Колкото популярността на страниците, които сочат към определена, е по-голяма, толкова и PR на тази страница е по-голям.


Page Rank Checker

[Изглед](#) [История](#) [Отметки](#) [Инструменти](#) [Помощ](#)

er - Check You X +

https://dnschecker.org/pagerank.php

 [Get Firefox Addon](#) [Donate](#) [f](#) [t](#)

[Home](#) [Flush DNS](#) [DNS Servers](#) [Reverse DNS Lookup](#) [All Tools](#) 

Your IP : 90.154.219.53


[DNS Checker](#) / [Page Rank Checker](#)

Pagerank Checker Tool

Pagerank checker tool enables the webmasters to lookup page rank of any website. Officially the Google PageRank service has been closed, but this tool tells you the last detected PageRank (if any) of a given website.

Enter any valid URL to check Page Rank.

Enter Domain to Check Page Rank:



Check Page Rank

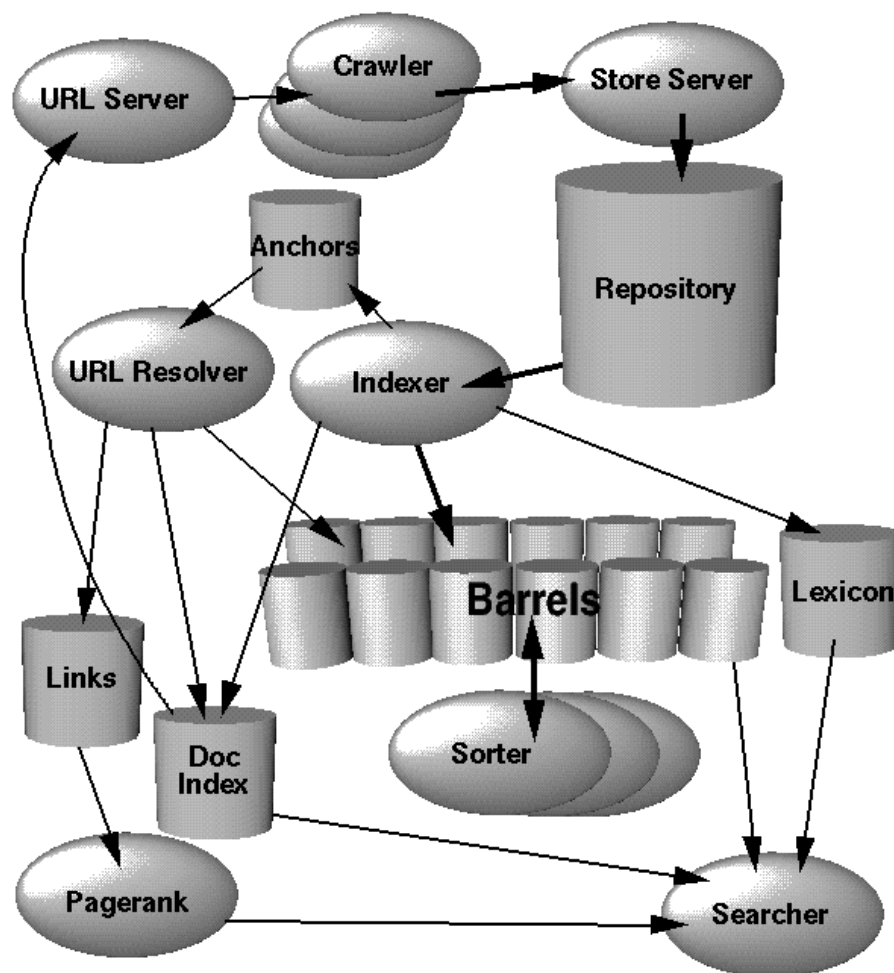
who.int

Page Rank

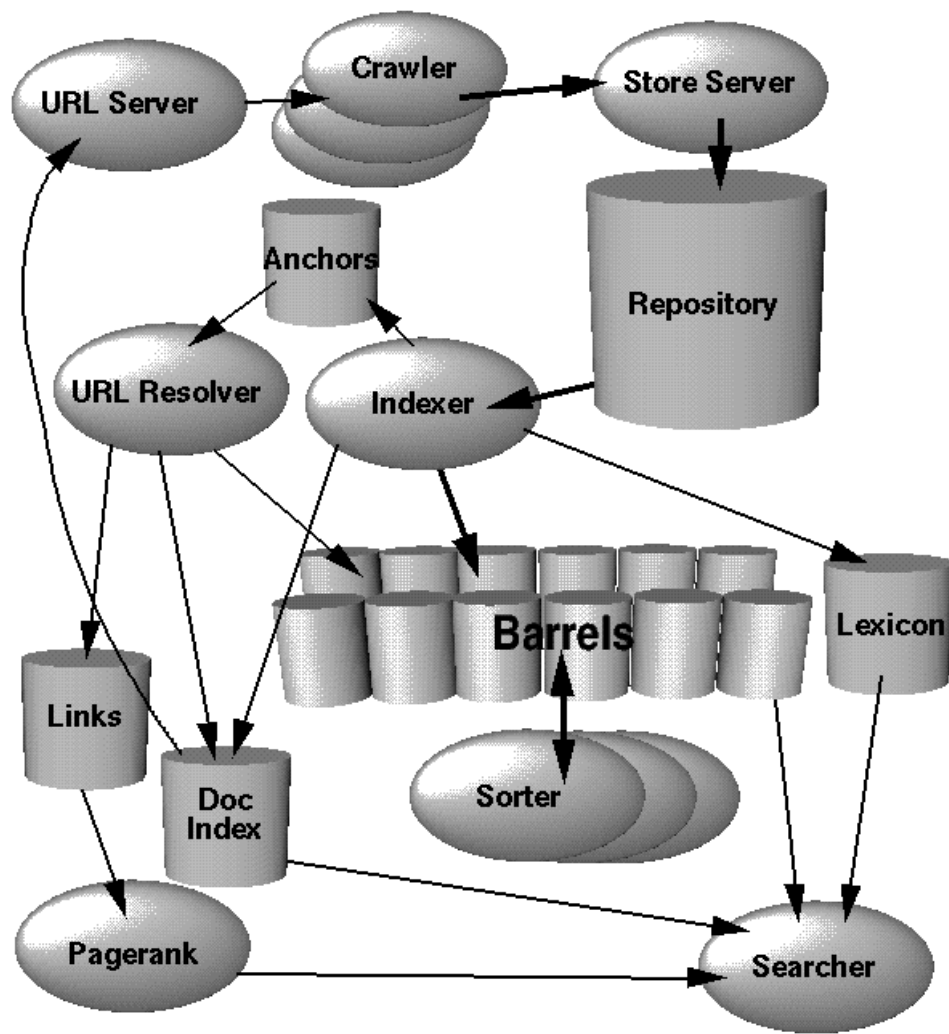
7 / 10

Интернет сървъри и технологии, Христо Вълчанов

Архитектура на Google

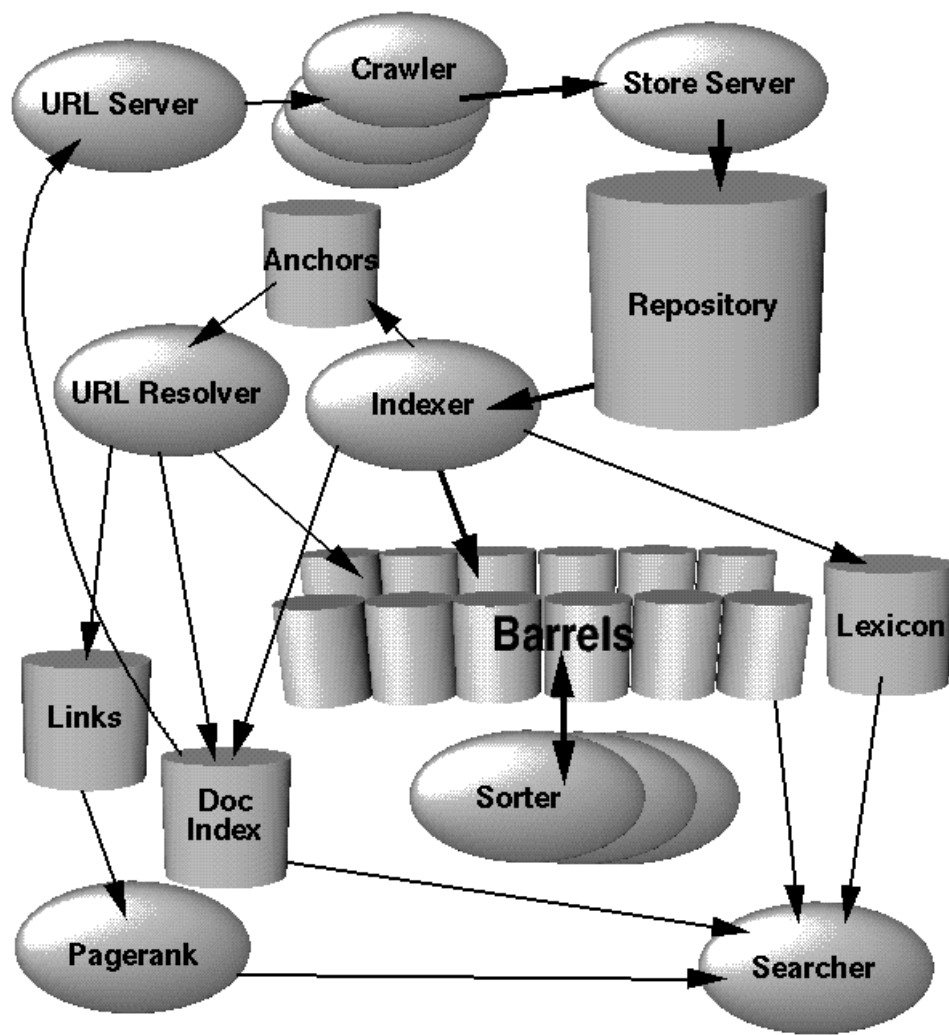


Компоненти - servers



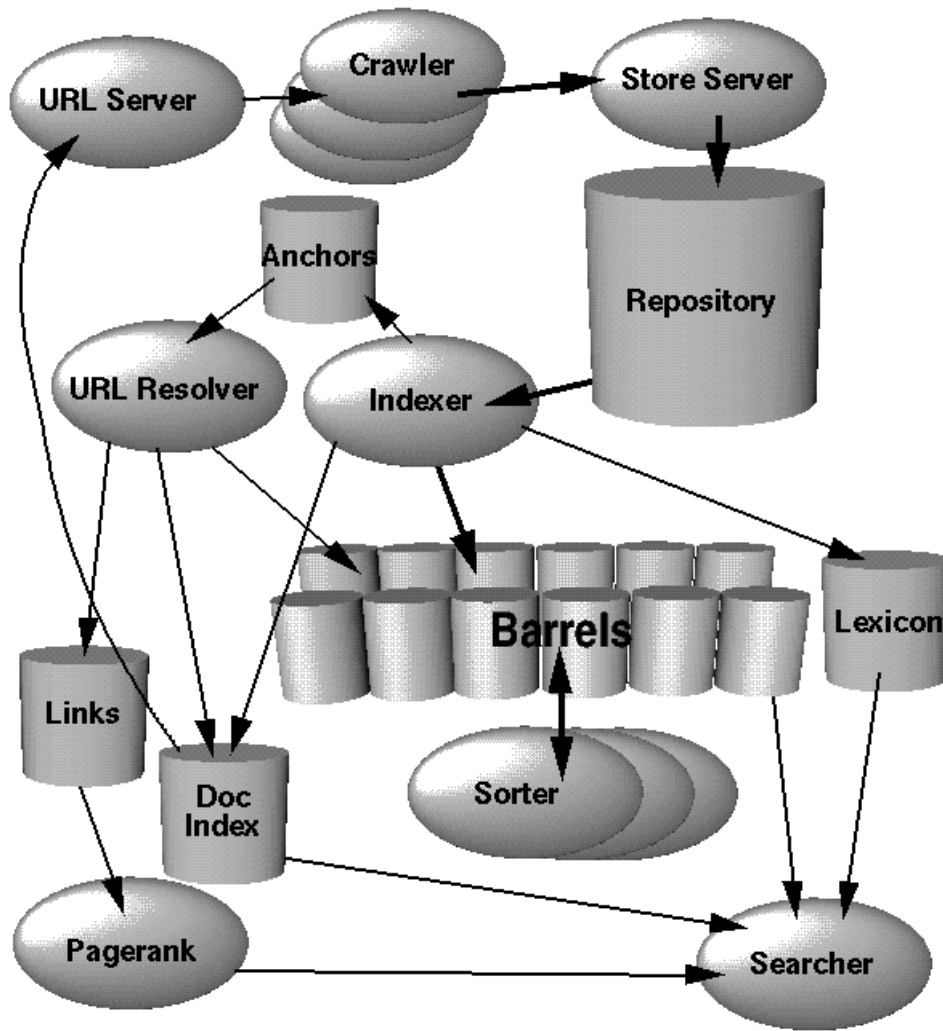
- URL server
- Store server

Компоненти - Repository



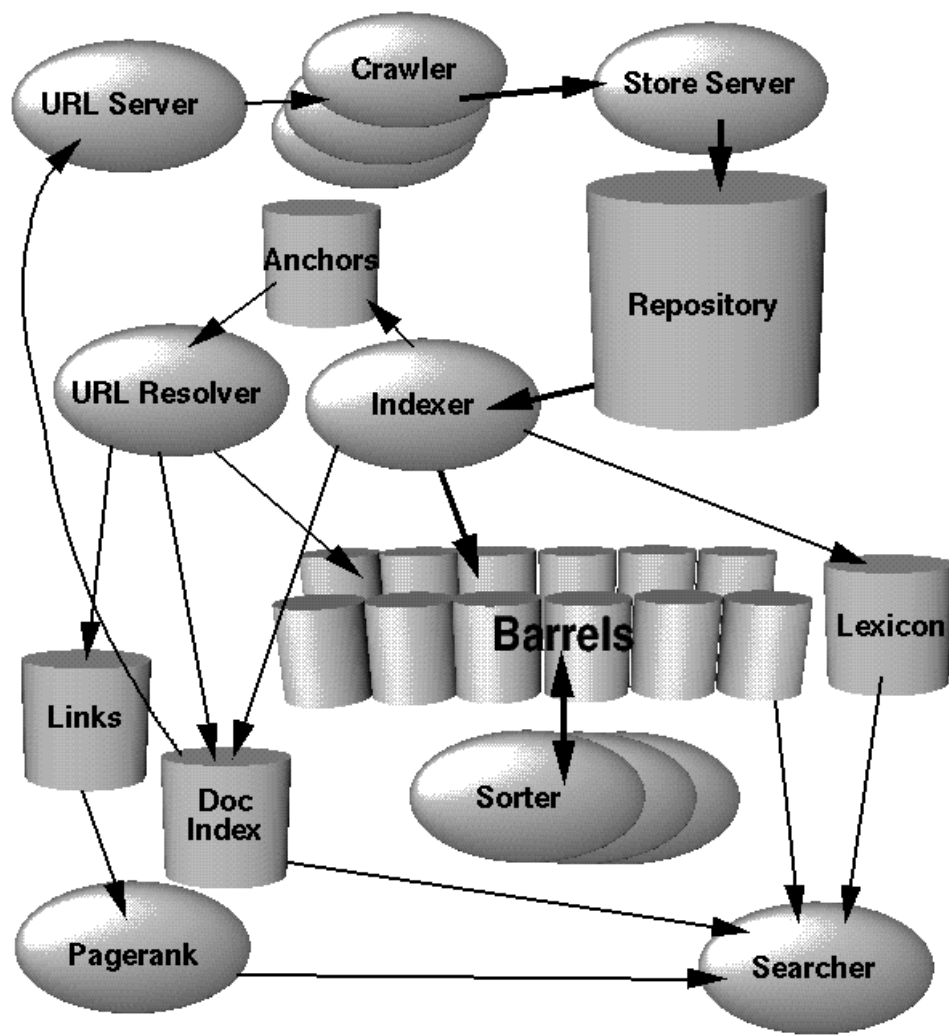
- Съдържа пълен HTML на всяка страница;
- Компресиране с zlib;
- Документите се съхраняват един след друг (docID, Len, URL, page)

Компоненти - Indexer



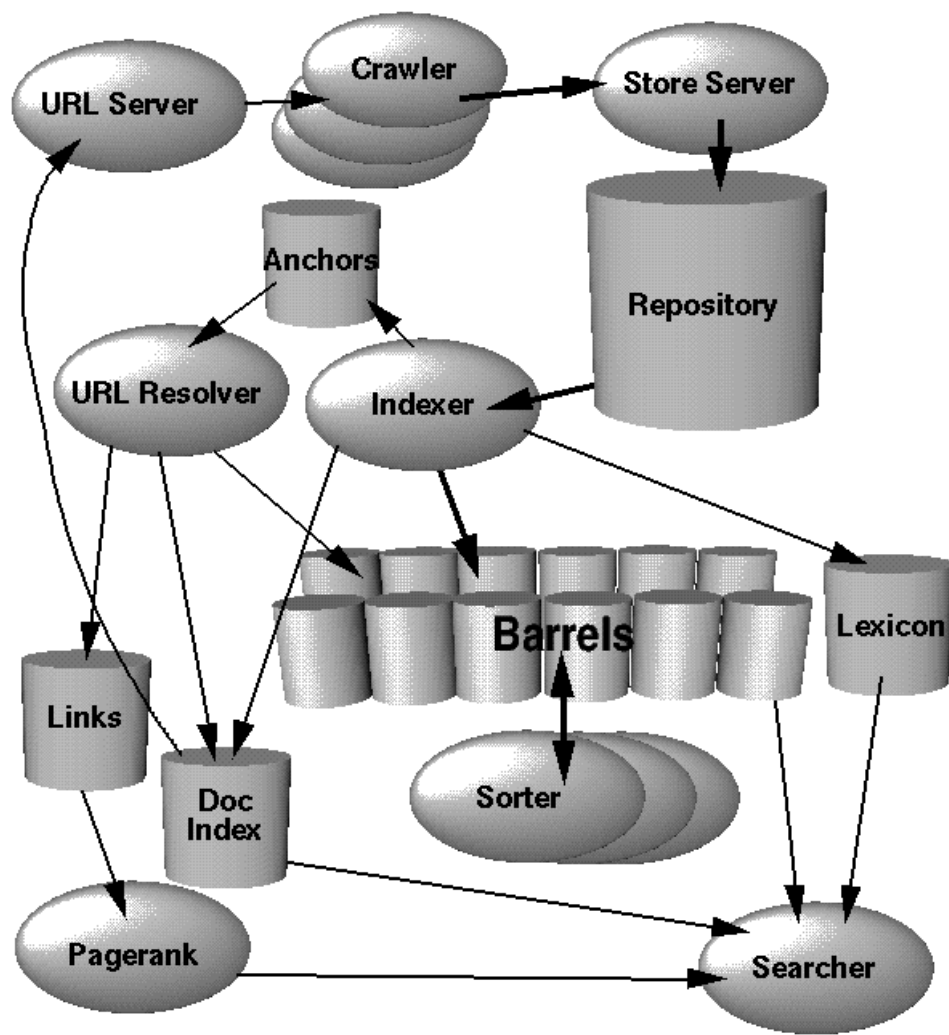
- Декомпресира документите.
- Всеки документ се конвертира в множество думи – *hits*.
- Hit – дума, позиция, размер шрифт, главни букви.
- Разпределят се в частично сортирани индекси – *barrels*.
- Парсва линковете в *anchor file*.

Компоненти – URL resolver



- Конвертира относителните URL в абсолютни;
- Поставя текста от котвата в индекс с docID, сочен от нея;
- Генерира база от линкове (линк-docID).

Компоненти – Sorter



- Сортира barrels по docID.
- Генерира инвертен индекс.
- Обновява лексикона.