

Introduction

IN₄₃₂₅ – Information Retrieval, 2012

Organizational matters

Course overview

- 2 lectures a week for 7 weeks
- 4 programming assignments
 - **Groups** of 2-3 students (find your own, email me your group name and members)
 - Assignments are handed out on Wednesdays; they are due the first following Wednesday with a class
- Final assignment
 - **Individual** work
 - Write a paper about an IR research topic
 - Hand-in intermediate results (topic, outline)

Assignments I-4

- Group work
- The Wednesday lecture will provide you with the background knowledge to do these assignments

Final assignment

- Individually
- Write an **8-page paper** about an information retrieval research topic of your choice
 - Can be a lecture topic
 - Start thinking about a potential topic early on in the course
- Intermediate deadlines will provide you with an opportunity to develop the topic and paper outline
 - Handing in intermediate results is **voluntary but strongly encouraged**
 - I will provide feedback until the middle of April

Assignments timeline

	Handed out	Due
Assignment 0	8/2/2012	----
Assignment 1	15/2/2012	22/2/2012
Assignment 2	22/2/2012	7/3/2012
Assignment 3	7/3/2012	14/3/2012
Assignment 4	14/3/2012	21/3/2012
Assignment 5	21/3/2012	First week of Q4
Assignment 5a	21/3/2012	28/3/2012
Assignment 5b	28/3/2012	15/4/2012

Grading

- Course grade: $0.5 * [A1/2/3/4] + 0.5 * A5$
- Assignments 0, 5a and 5b are voluntary and do not count towards your grade
- To pass
 - an overall grade of 6 or higher is required
 - AND all required assignments are handed in on time
 - AND all required assignments are passed with a score of 6 or higher

Questions?

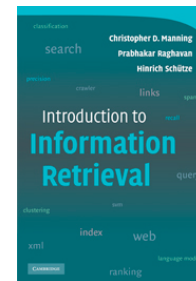
- Are always welcome!
- Always per email c.hauff@tudelft.nl
 - Emails with **[IN4325]** in the subject line are read first!!
- Questions & answers may be posted on blackboard if they are relevant for others as well!

Reading material (lectures)

Additional papers will be announced on blackboard if necessary.

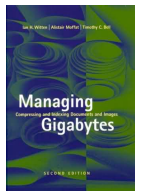
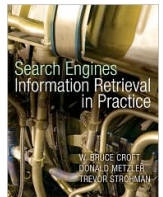
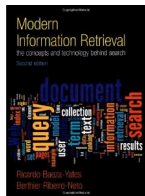
- *Introduction to Information Retrieval* by Manning, Raghavan and Schütze, University Press, 2008.

- Book used in this course
- Available online: <http://nlp.stanford.edu/IR-book/>



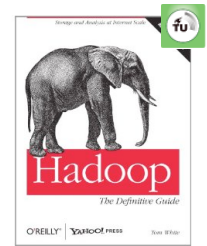
- Additional suggestions

- *Modern Information Retrieval* by Baeza-Yates and Ribeiro-Neto, Addison-Wesley Professional; 2nd edition, 2011.
- *Search Engines: Information Retrieval in Practice* by Croft, Metzler and Strohman, Addison-Wesley, 2009.
- *Managing Gigabytes: Compressing and Indexing Documents and Images* by Witten, Moffat and Bell, Morgan Kaufmann, 1999.



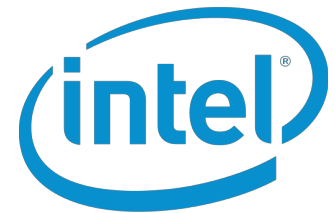
Reading material (assignments)

- The assignments are focused on MapReduce & Hadoop (more about this in the next lecture)
- Hadoop programming
 - *Hadoop: The Definitive Guide* by Tom White, O'Reilly Media, 2011.
- MapReduce algorithm design
 - *Data-Intensive Text Processing with MapReduce* by Lin, Dyer and Hirst, Morgan and Claypool Publishers, 2010
 - Available online: <http://www.umiacs.umd.edu/~jimmylin/book.html>



Information retrieval

Information retrieval in industry



Information retrieval definitions

- “The goal of a machine method of information retrieval is purely and simply that of being able to **find** and to **recover** at will information stored in a **collection of documents**. [...] It is oriented completely towards actual use of the information, and to the **convenience of the user**.”

Calvin N. Mooers, Scientific information retrieval systems for machine operation; case studies in design. XIIth International Congress of Pure and Applied Chemistry, 1951

Information retrieval definitions

- “The goal of a machine method of information retrieval is purely and simply that of being able to **find** and to **recover** at will information stored in a **collection of documents**. [...] It is oriented completely towards actual use of the information, and to the **convenience of the user**.”

Calvin N. Mooers, Scientific information retrieval systems for machine operation; case studies in design. XIIth International Congress of Pure and Applied Chemistry, 1951

- “An information retrieval system does not inform (i.e. change the knowledge of) the user on the subject of his inquiry. It merely informs on the **existence (or non-existence)** and **whereabouts** of documents relating to his request.”

F.W. Lancaster. Information Retrieval Systems: Characteristics, Testing and Evaluation, Wiley, New York, 1968.

Information retrieval definitions II

- “Information retrieval studies the retrieval of information (not data) from a collection of written documents. The retrieved documents aim at satisfying a **user information need** usually expressed in **natural language**.”

Baeza-Yates and Ribeiro-Neto. Modern Information Retrieval, 1999

Information retrieval definitions II

- “Information retrieval studies the retrieval of information (not data) from a collection of written documents. The retrieved documents aim at satisfying a **user information need** usually expressed in **natural language**.”

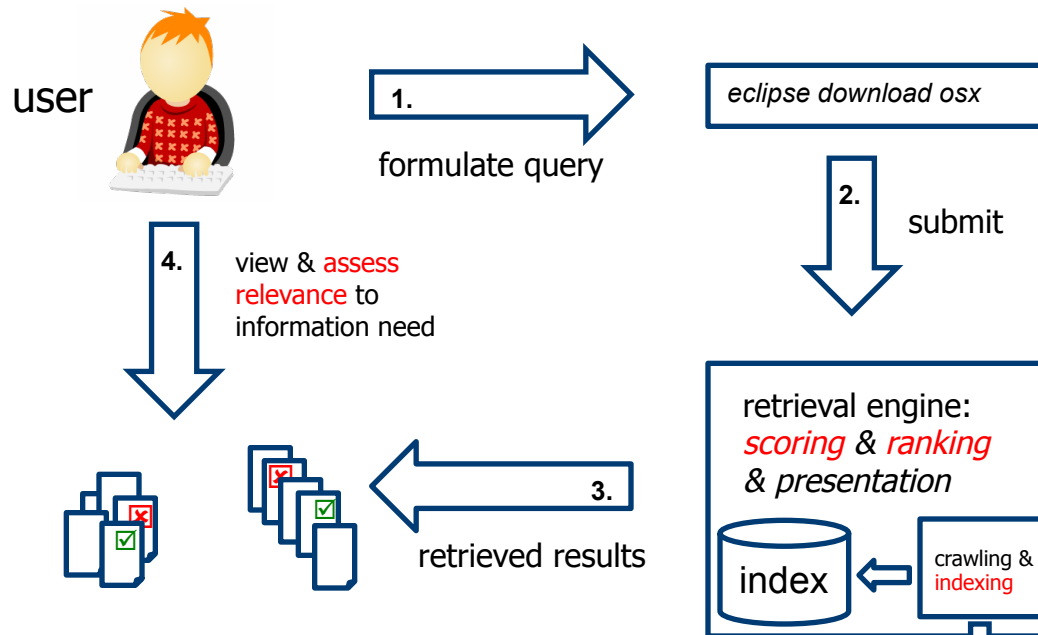
Baeza-Yates and Ribeiro-Neto. Modern Information Retrieval, 1999

- “IR has traditionally been concerned primarily with the process of representation of texts and queries, and a comparison of these representations. [...] it is becoming increasingly evident that IR is an inherently **interactive process** [...] This means, in particular, that supporting and taking advantage of the interaction of the user with the other components of the IR system is crucial for effective IR system design.”

N.J. Belkin and C. Cooll, Cases, Scripts, and Information-Seeking Strategies: On the Design of Interactive Information Retrieval Systems. Expert Systems With Applications, Vol. 9, No. 3, 1995

This course in short (red)

information need: *I am supposed to use Eclipse for the assignments. Where can I download a version for Mac OS X?*



Chapter 13. The Hundred Days.

M. Noiret was a true prophet, and things progressed rapidly, as he had predicted. Every one knows the history of the famous return from Elba, a return which was unprecedented in the past, and will probably remain without a counterpart in the future.

Louis XVIII. made but a faint attempt to parry this unexpected blow; the monarchy he had scarcely reconstructed tottered on its precarious foundation, and at a sign from the emperor the incongruous structure of ancient prejudices and new ideas fell to the ground. Villéfort, therefore, gained nothing save the king's gratitude (which was rather likely to injure him at the present time) and the cross of the Legion of Honor, which he had the pride not to wear, although M. de Biscan had duly forwarded the investiture.

Napoleon would, doubtless, have deprived Villéfort of his office had it not been for Noiret, who was all powerful at court, and thus the Girondin of '93 and the Senator of 1806 protected him who so lately had been his protector. All Villéfort's influence barely enabled him to stifle the secret Dantes that so nearly divulged. The king's procurator alone was deprived of his office, being suspected of royalism.

However, scarcely was the imperial power established—that is, scarcely had the emperor re-entered the Tuilleries and begun to issue orders from the closet into which we have introduced our readers,—he found on the table there Louis XVIII's half-filled snuff-box,—scarcely had this occurred when Maresilles began, in spite of the authorities, to rekindle the flames of civil war, always smouldering in the south, and it required but little to excite the populace to acts of far greater violence than the shouts and insults with which they assailed the royalists wherever they ventured abroad.

Owing to this change, the worthy shipowner became at that moment—we will not say all powerful, because Morel was a prudent and rather a timid man, so much so, that many of the most zealous partisans of Bonaparte accused him of "moderation"—but sufficiently influential to make a demand in favor of Dantes.

Villéfort retained his place, but his marriage was put off until a more favorable opportunity. If the emperor remained on the throne, Gérard required a different alliance to aid his career: if Louis XVIII. returned, the influence of M. de Saint-Méran, like his own, could be vastly increased, and the marriage he still more suitable. The deputy-procurator was, therefore, the first magistrate of Marseilles, when one morning his door opened, and M. Morel was announced.

Any one else would have hastened to receive him; but Villéfort was a man of ability, and he knew this would be a sign of weakness. He made Morel wait in the ante-chamber, although he had no one with him, for the simple reason that the king's procurator always makes every one wait, and after passing a quarter of an hour in reading the papers, he ordered M. Morel to be admitted.

Morel expected Villéfort would be dejected; he found him as he had found him six weeks before, calm, firm, and full of that glacial politeness, that most insurmountable barrier which separates the well-bred from the vulgar man.

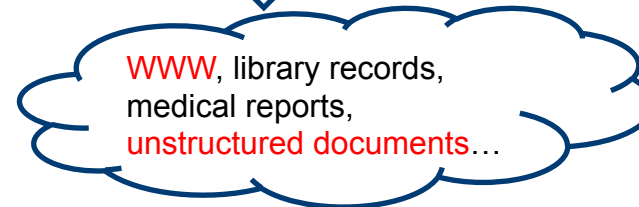
Example of unstructured text.

Not covered [for now]:

- result presentation (clustering vs. ranked list, snippets, interfaces)
- crawling
- cross-lingual IR
-

Related fields:

- Database management (more structured data)
- Library & information science (categorization)
- Semantic web (logic, reasoning)
- Natural language processing
- Machine learning (text classification & mining)



Why so complicated?

- Searching for the lines in the book *Count of Monte Christo* that contain the terms *Dantes* AND *prison* but NOT *Albert*
- Naïve solution
 - Grep all lines that contain *Dantes*, then grep those containing *prison* and finally strip out lines containing *Albert*

```
more countOfMonteChristo.txt|grep Dantes|grep prison|grep -v Villefort
```

- Problems
 - Proximity operations not easy to implement (e.g. *Dantes* within max. 3 terms of *prison*)
 - *Set of matching results (yes/no decision)*
 - What about approximate/semantic matches (Edmond instead of Dantes, cell instead of prison. etc.)

Why so complicated?

- Searching for the lines in the book *Count of Monte Christo* that contain the terms *Dantes* AND *prison* but NOT *Albert*
- Naïve solution
 - Grep all lines that contain *Dantes*, then grep those containing *prison* and finally strip out lines containing *Albert*

```
more countOfMonteChristo.txt|grep Dantes|grep prison|grep -v Villefort
```

Elaborate queries require the user to anticipate possibly used terms:

(*Edmond* OR *Dantes* OR *Monte-Cristo*) AND
(*prison* OR *cell* OR *imprisoned*) NOT *Albert*

- What about approximate/semantic matches (Edmond instead of Dantes, cell instead of prison. etc.)

Why so complicated? II

- What about using a *term-document* matrix?
 - Here: a line is a document

	L1	L2	L3	L4	L5	L6	L7	L8	L9
Edmond	0	0	0	1	1	0	1	0	1
Dantes	1	0	1	0	1	0	1	0	0
Monte-Cristo	0	1	0	0	0	1	0	0	0
prison	1	0	0	1	0	0	1	1	0
cell	0	1	1	0	0	0	1	0	0
imprisoned	1	0	0	0	0	0	1	0	0
Albert	1	1	0	0	1	0	0	1	0

1 if *Edmond* occurs in line L9; 0 otherwise

Only feasible for extremely small corpora. The matrix gets too large too quickly. Still no ranked retrieval.

Dantes AND *prison* \Rightarrow bitwise AND

1	0	1	0	1	0	1	0	0
1	0	0	1	0	0	1	1	0

1	0	0	0	0	0	1	0	0

\Rightarrow L1 L7

Dantes NOT *Albert* \Rightarrow bitwise AND complement

1	0	1	0	1	0	1	0	0
0	0	1	1	0	1	1	0	1

0	0	1	0	0	0	1	0	0

\Rightarrow L3 L7

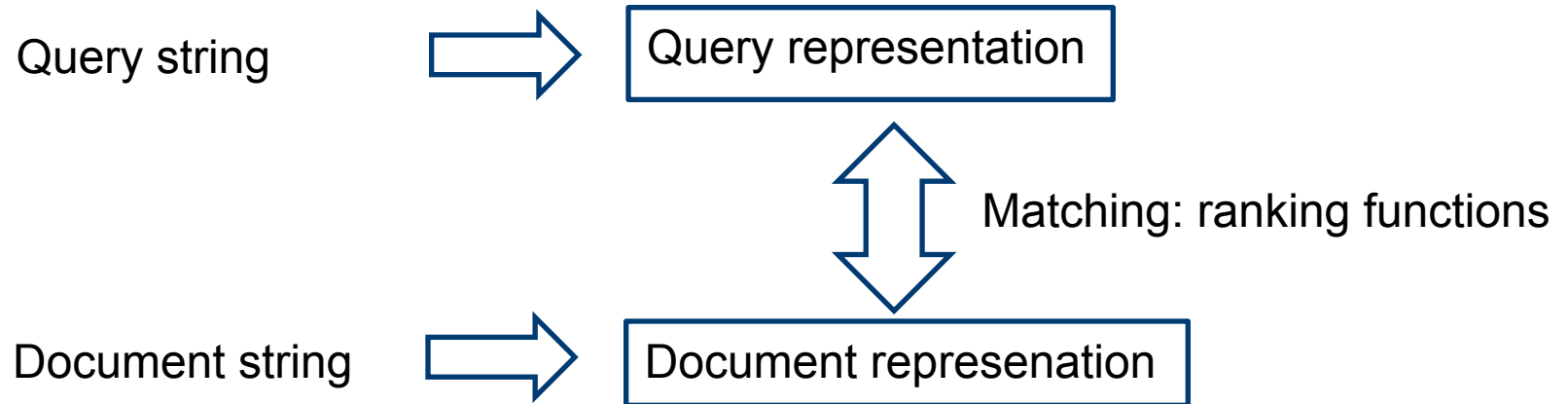
Too many vs. too few

~50,000 lines of text in the book

Boolean queries	#Lines retrieved
<i>Edmond</i> OR <i>Dantes</i> OR <i>Monte-Cristo</i>	1995
<i>Edmond</i> OR <i>Dantes</i> OR <i>Monte-Cristo</i> NOT <i>Albert</i>	1976
<i>prison</i> OR <i>cell</i> OR <i>imprisoned</i>	587
<i>prison</i>	240
(<i>Edmond</i> OR <i>Dantes</i> OR <i>Monte-Cristo</i>) AND (<i>prison</i> OR <i>cell</i> OR <i>imprisoned</i>)	22
<i>Edmond</i> AND <i>prison</i>	3

Why so complicated? III

Since simple string matching is not enough, we model ...



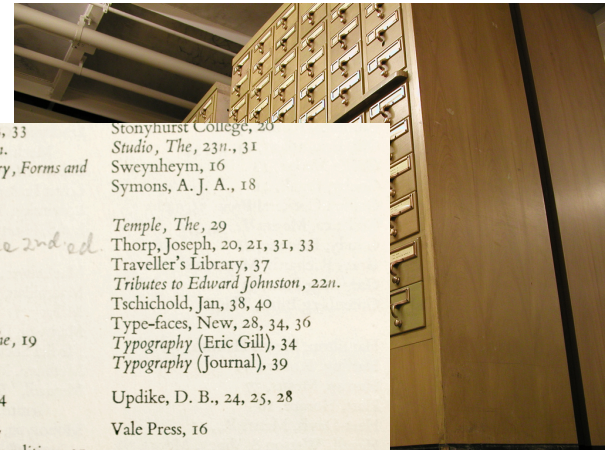
A little bit of IR history

Information retrieval roots

Library classification models and card catalogues (1950's)



Library of Congress Classification



Printing of Books, *The*, 32, 33
Printing Review, 2011, 2111.
Printing Types, *their History, Forms and Use*, 25
Processes, New, 17
Psalter, Chantilly, 12
Ravilious, Eric, 20
Ricketts, Charles, 16, 18
Ridder, Vivian, 18
Rogers, Bruce, 24, 28, 38
Roots of the Mountains, *The*, 19
Royal College of Art, 23
Ruskin, John, 9, 11
Rutherford, Albert, 32, 34
Secker, Messrs Martin, 37
Shakespeare, One-volume edition, 31
Shakespeare Head Press, 30, 31
Stonyhurst College, 20
Studio, *The*, 2311, 31
Sweynheym, 16
Symons, A. J. A., 18
Temple, *The*, 29
Thorpe, Joseph, 20, 21, 31, 33
Traveller's Library, 37
Tributes to Edward Johnston, 2211.
Tschichold, Jan, 38, 40
Type-faces, New, 28, 34, 36
Typography (Eric Gill), 34
Typography (Journal), 39
Updike, D. B., 24, 25, 28
Vale Press, 16
Warde, Beatrice, 21, 28
Wardrop, James, 2211, 23

Book index creation: a very laborious task if done manually.

NY Times Print Index

The goal of classification schemes (many exist) is to remove the ambiguity of natural language.

Documents are indexed/described by human annotators.

Universal decimal classification




expand all collapse all	
0 SCIENCE AND KNOWLEDGE. ORGANIZATION. COMPUTER SCIENCE. INFO	
00 Prolegomena. Fundamentals of knowledge and culture. Propaedeutics	
001 Science and knowledge in general. Organization of intellectual work	
001.1 Concepts of science and knowledge	
001.18 Future of knowledge	
001.32 Learned, scientific societies. Academies	
001.8 Methodology	
001.89 Organization of science and scientific work	
001.9 Dissemination of ideas	
002 Documentation. Books. Writings. Authorship	
003 Writing systems and scripts	
004 Computer science and technology. Computing. Data processing	
004.01/.08 Special auxiliary numbers for computing	
004.2 Computer architecture	
004.3 Computer hardware	
004.4 Software	
004.5 Human-computer interaction. Man-machine interface. User	
004.6 Data	
004.7 Computer communication. Computer networks	
004.8 Artificial intelligence	
004.9 Application-oriented computer-based techniques	
005 Management	
006 Standardization of products, operations, weights, measures and t	
007 Activity and organizing. Communication and control theory gener	
008 Civilization. Culture. Progress	
01 Bibliography and bibliographies. Catalogues	
02 Librarianship	

004	Computer science and technology. Computing. Data processing
004.2	Computer architecture
004.22	Data representation
004.23	Instruction set architecture
004.25	Memory system
004.27	Advanced architectures. Non-Von Neumann architectures

<http://www.udcc.org>

Vannevar Bush (1890-1974)

Early visions of the Web in 1945

"Consider a future device for individual use, which is a sort of mechanized private file and library. It needs a name, and, to coin one at random, "**memex**" will do. [...] a device in which an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with  exceeding speed and flexibility. **It is an enlarged intimate supplement to his memory.** [...] It affords an immediate step, however, to associative indexing, the basic idea of which is a provision whereby **any item may be caused at will to select immediately and automatically another.**  Thereafter, at any time, when one of these items is in view, the other can be instantly recalled merely by tapping a button below the corresponding code space. [...] It is exactly as though the physical items had been gathered together from widely separated sources and bound together to form a new book. It is more than this, for **any item can be joined into numerous trails.** [...] **Wholly new forms of**  **encyclopedias will appear, ready made with a mesh of associative trails running through them.** [...] There is a new profession of **trail blazers**, those who find delight in the task of establishing useful trails through the enormous mass of the common record."



Medlars

Launched in 1964

- Medical Literature Analysis and Retrieval Service: computer-based retrieval system for medical research literature

- Searching based on human-assigned indexing terms
- Boolean logic queries (*cancer AND mouse*)
- Search process:

In general: librarians, paralegals, patent officers, etc.

1. User (e.g. medical researcher) corresponds with **expert searcher** via mail, phone or face-to-face [**delegated search vs. end-user search**]
2. Expert searcher formulates query
3. Batch run of queries overnight on dedicated machine
4. Result printouts posted to the user

Querying the system could take **days or weeks**.

newcastle
newcastle disease
newcastle disease virus
newcastle-ottawa
newcastle-ottawa scale
newcastle university
newcastle disease virus cancer
newcastle disease vaccine
newcastle vaccine

Turn off

Search

Help

Display Settings: Abstract

Indian J Med Res. 2009 Nov;130(5):507-13.

Newcastle disease virus as an on

Ravindra PV, Tiwari AK, Sharma B, Chauhan RS.

Molecular Biology Laboratory, Division of Animal Biotechno

Abstract

Cancer is a major cause of deaths in humans. The current chemo- and radiotherapies have provided treatment of cancer termed, oncolytic virotherapy has recently emerged. Newcastle disease virus (NDV) is one such virus with an inherent oncolytic property. NDV causes a highly infectious disease in poultry worldwide. In humans it is reported to have oncolytic and immuno-stimulatory effects. It specifically replicates in tumour cells while sparing normal cells and cause oncolysis. For many years different strains of the NDV have been investigated for treatment of various human cancers. Recent advances in reverse genetics provided investigators the tools to produce recombinant NDV with improved oncolytic property.

PMID: 20090097 [PubMed - indexed for MEDLINE] Free full text

Publication Types, MeSH Terms

Publication Types

Research Support, Non-U.S. Gov't

Review

MeSH Terms

Animals

Apoptosis

Humans

Neoplasms/pathology

Neoplasms/therapy*

Newcastle disease virus/genetics

Newcastle disease virus/physiology*

Oncolytic Virotherapy/methods*

Oncolytic Viruses/genetics

Oncolytic Viruses/physiology

LinkOut - more resources

Manually added
by experts!
Even today.

Send to:

IJMR
Free Full Text

Related citations

Review [Progress in using Newcastle disease virus for [Sheng Wu Gong Cheng Xue Bao. 2010]

Type I interferon-sensitive recombinant newcastle disease virus for oncolyti [J Virol. 2010]

Generation of a recombinant oncolytic Newcastle disease virus and expression c [Gene Ther. 2008]

Analysis of three properties of Newcastle disease virus for fighting cancer: [Methods Mol Biol. 2012]

Review Newcastle disease virus (NDV): brief history of its oncolytic strains. [J Clin Virol. 2000]

See reviews...

See all...

Cited by 1 PubMed Central article

Cytolytic replication of echoviruses in colon cancer cell lines. [Virol J. 2011]

Related information

Related Citations

Cited in PMC

Recent activity

Turn Off Clear

Newcastle disease virus as an oncolytic agent. PubMed

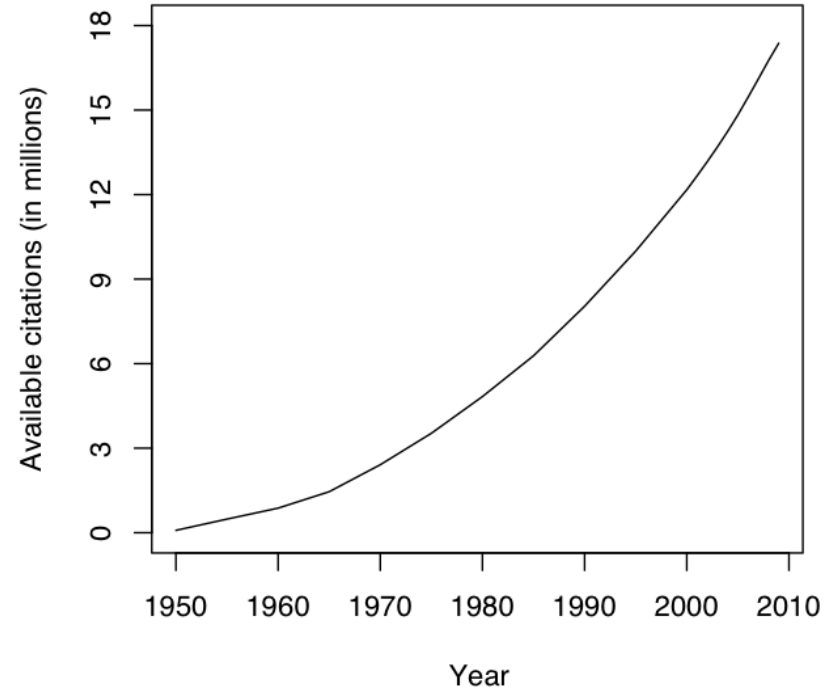
Clinical Evidence for the Regression of Liver Fibrosis. PubMed

Non-drug therapies for lower limb muscle cramps. PubMed

Variability among the neuraminidase, non-structural 1 and PB1-F2 proteins i PubMed

They like to publish ...

- MEDLINE: part of PubMed
 - Bibliographic database
 - Exponential growth
 - 2010 statistics
 - 18 million citations
 - 5,516 journals
 - 700,000 additions



Slide provided by Dolf Trieschnigg, University of Twente.

SMART project

1960s-1990s

- Gerard Salton (1927-1995)
 - “the man most responsible for the establishment, survival, and recognition of Information Retrieval” [8]
- Developed pioneering solutions that are still in use today
 - Fully automatic indexing of document texts
 - Scoring functions to determine how well the query and document match
 - The concept of result ranking (vs. result set)
 - Relevance feedback (exploit known relevant documents)
- Scoring & ranking were added to commercial systems in the late 1980s and became ubiquitous with web search engines in the 1990s

SMART project

1960s-1990s

Computing power & costs: in **1973** one 'run' on the **1400-document** Cranfield collection took **11.2** minutes in processing time and it cost **\$86.22** [1].

- Gerard Salton (1927-1995)
 - “the man most responsible for the establishment, survival, and recognition of Information Retrieval” [8]
- Developed pioneering solutions that are still in use today
 - Fully automatic indexing of document texts
 - Scoring functions to determine how well the query and document match
 - The concept of result ranking (vs. result set)
 - Relevance feedback (exploit known relevant documents)
- Scoring & ranking were used in the 1980s and became ubiquitous in the 1990s

“... it is nice to be able, in England, to search a file in California of **five hundred thousand documents** using an arbitrarily complex search specification, and get answers back in **six seconds**.”
(Karen Spärck Jones, **1979**)

Okapi, BM₂₅ and Okapi BM₂₅

1990s-2000s

- Stephen Robertson & collaborators (London City University)
- Okapi is a retrieval system with ranking function BM₂₅
- Focus on weighting terms, exploiting user feedback and user-system interaction

GAT-INQUIRY

Type your query below and click on "Search"

how has affirmative-action affected the construction-industry construction projects public works

Search Clear all

Query results

0 docs saved. 0 docs relevant.

- Save? Rel? 1. AP: Court Refuses to Kill
- Save? Rel? 2. AP: Court To Decide Affirmative Action
- Save? Rel? 3. AP: Tennessee Families Challenging 'G
- Save? Rel? 4. Sn Js: TOO FEW WOMEN SUCCEED IN
- Save? Rel? 5. Fed Register: No title for this document
- Save? Rel? 6. AP: Officials Express Regret Over Reje
- Save? Rel? 7. W St J: Labor Letter: A Special News
- Save? Rel? 8. Fed Register: No title for this document
- Save? Rel? 9. AP: Court Refuses to Kill Florida Affirm
- Save? Rel? 10. Fed Register: No title for this document
- Save? Rel? 11. W St J: International: U.S., Japan S
- Save? Rel? 12. W St J: Industry Focus: U.S. Contra
- Save? Rel? 13. Fed Register: No title for this document
- Save? Rel? 14. AP: High Court Upholds Limits on Loc
- Save? Rel? 15. Fed Register: No title for this document
- Save? Rel? 16. Fed Register: No title for this document
- Save? Rel? 17. Fed Register: No title for this document
- Save? Rel? 18. Fed Register: No title for this document
- Save? Rel? 19. Fed Register: No title for this document
- Save? Rel? 20. AP: Supreme Court Examines Affirma
- Save? Rel? 21. Fed Register: No title for this document
- Save? Rel? 22. AP: Delay Urged In Action Aimed At O
- Save? Rel? 23. AP: Government Hails Construction A
- Save? Rel? 24. AP: Construction Executives Call For
- Save? Rel? 25. AP: Kennedy Could Be Key Vote on M
- Save? Rel? 26. AP: Kennedy May Hold Decisive Vote
- Save? Rel? 27. W St J: Hong Kong Frictions Grow Ov
- Save? Rel? 28. W St J: Law: Court Rejects Schol
- Save? Rel? 29. Fed Register: No title for this document
- Save? Rel? 30. W St J: Economy: Fed Reports Wea

Current logfile: [/homes/mg/Trec6Logs/Okapi/t326.mjg.0]

To add terms to the query type: (a) one or more words, or (b) one phrase ending in a + sign, then press return

Working Query		Document Hitlist
35	2 : ferry sinking (B)	36: FT943-3397 [713] 1/1 page
265	2 : loss of life (B)	FT 14 SEP 94 / UK Company News: United Friendly ahead sharply to Pounds 13.6m A turnaround in the general insurance business underpinned a sharp rise at United Friendly, wher.....
2525	2 : disaster	loss (2) life (2)
44641	3 : operators	37: FT923-2285 [710] 1/1 page
31463	3 : loss	FT 18 SEP 92 / UK Company News: Reorganisation moves help put L&G Pounds 74m back in black LEGAL & GENERAL, the life assurance group, reported a turnaround to pre-tax profits.....
16963	2 : life	life (4) loss (7)
		38: FT941-14279 [709] 1/1 page
		FT 21 JAN 94 / Clinton gives more help to quake victims President Bill Clinton yesterday gave California a Dollars 100m (Pounds 67.5m) advance for earthquake repairs and ann.....
		disaster (1) loss of life (2)
		39: FT931-373 [708] 1/1 page
		[F] 1000 FT943-178
		FT 30 SEP 94 / Leading Article: Defying the cruel sea Ferries are among the safest vessels afloat. But, as the tragic sinking of the Estonia with the loss of more than 800 l.....
		[F] 874 FT943-312
		FT 30 SEP 94 / Ferries in six 'near accidents': Finland and Sweden order checks after Estonia sinking STOCKHOLM, TALINN Sweden's government disclosed yesterday that six rece.....
		[F] 715 FT934-8043
		FT 17 NOV 93 / International Company News: Uni Storebrand back in black at nine months OSLO UNI Storebrand, Norway's biggest insurance group, yesterday reported nine-month p.....

Clear Current Query
 Clear Relevance Feedback
 Set Working Query Size
 Cancel Menu ^C

Search Database
Query Options
Exit Okapi

A. Veerasamy, Interacti

M.M. Beaulieu and M.J. Gatford, Interactive Okapi at TREC-6, 1997

TREC

Text REtrieval Conference (1992-*)

- Conducted by the US National Institute of Standards and Technology, co-sponsored by DARPA
- Several “tracks” per year (a good way to learn about current work)

New queries over the same data.

The same queries over new data.

Non-English texts.

OCR simulation

Documents of languages A & B are searched with queries of language C.

Ad-hoc retrieval task (1992)

Routing task (1992)

Interactive track (1995)

Multilingual track (1995)

Database merging track (1995)

Confusion track (1995)

Cross-Language track (1997)

Spoken document track (1997)

Question Answering track (1999)

Web track (1999)

Search for organizations, people, products.

Video track (2001)

Novelty track (2002)

Genomics track (2003)

Terabyte track (2004)

Enterprise track (2005)

Spam track (2005)

Blog track (2006)

Legal track (2006)

Million query track (2007)

Chemical IR track (2009)

Entity track (2009)

Microblog track (2011)

Now split from TREC: TRECVID.

Find the relevant and novel docs.

Enterprise data (intranet, email).

Business records.

Search for chemical patents and research articles.

Benchmarks are very important to IR

The largest ongoing benchmarks apart from TREC/TRECVID

- CLEF
 - Conference and Labs of the Evaluation Forum
 - <http://www.clef-initiative.eu/>
- MediaEval
 - Benchmarking Initiative for Multimedia Evaluation
 - <http://www.multimediaeval.org/>
- NTCIR
 - NII Test Collection for IR Systems
 - <http://research.nii.ac.jp/ntcir/index-en.html>
- FIRE
 - Forum for Information Retrieval Evaluation
 - <http://www.isical.ac.in/~clia/>



NTCIR



MediaEval Benchmark

Current systems and technologies

Current systems and technologies I

Auto-completion

Google

Search

Everything

Images

Maps

Videos

News

Shopping

More

All results

Sites with images

More search tools

dam square amsterdam
dam square amsterdam
dam square amsterdam hotels
dam square amsterdam address
dam square amsterdam the netherlands

Dam Square
www.dutchsouvenirs.com/
3 Google reviews - Write a review

Dam Square 17 1012 JS Amsterdam, Netherlands
020 6203432

Dam Square - Wikipedia, the free encyclopedia
en.wikipedia.org/wiki/Dam_Square
Dam Square, or simply the Dam (Dutch: de Dam) is a town square in **Amsterdam**, the capital of the Netherlands. Its notable buildings and frequent events make it ...
Location and description - History - Dam Square shooting, 1945 - Present

Amsterdam.info - Dam Square Amsterdam hotels
www.amsterdam.info > Sights > Squares
The **Dam square** is the very central square of **Amsterdam**.

Dam Square, Amsterdam - Things to Do - VirtualTourist
www.virtualtourist.com > ... > Amsterdam > Things to Do
Dam Square reviews and photos from real travelers and locals in **Amsterdam**, Netherlands.

Dam Square - Amsterdam - Reviews of Dam Square - TripAdvisor
www.tripadvisor.com/Attraction_Review-g188590-d189381-Review...
★★★★☆ 84 reviews
Dam Square, Amsterdam: See 84 reviews, articles, and 110 photos of Dam Square, ranked No.87 on TripAdvisor among 458 attractions in Amsterdam.

Dam Square Amsterdam, Netherlands
www.360cities.net > The World > Europe > Netherlands > Amsterdam
10 Jan 2006 - 360° panoramic photography by Jook Leung | 360VR Images. Visit us to see more amazing panoramas from **Amsterdam** and thousands of ...

Dam Square - Directions

Query intent

Details

Transit: Dam

More reviews

virtualtourist.com (38)
flipkey.com (39)
tripadvisor.co.uk (40)
tripadvisor.com (11)

Where is the spam?

Ads adaptation (none shown)

Current systems and technologies II

The image shows a Google search interface for the query "flight amsterdam dublin". The search bar is at the top, with the Google logo on the left and a search button on the right. Below the search bar, the text "Search" is displayed in red, followed by "About 9,440,000 results (0.35 seconds)". A blue callout box points to the search results area, containing the text "~9 million results: ranking is of utmost importance".

On the left side of the search results, there is a vertical navigation menu with the following items: "Everything", "Images", "Maps", "Videos", "News", "Shopping", "More", and "Show search tools".

The main search results area is divided into two columns. The left column contains several search results, including:

- Amsterdam - Dublin? - Goedkoop Vliegen naar Dublin?**
www.cheaptickets.nl/Dublin
Boek Nu Online: **Amsterdam - Dublin!**
Vertrouwd boeken - Allerlaagste tarieven - Goedkope Citybreaks
- Greatest Offers on - Amsterdam-Dublin Flights | eDreams.com**
www.edreams.com/Amsterdam_Dublin
Book Now -Limited Seats Available!
► Early 2012 Flight Specials - Amsterdam-London from €20
- Flight Amsterdam Dublin - Op 1 website alle vliegtickets**
www.worldticketcenter.nl/AMS-DUB
Goedkoop en snel geregeld

The right column contains several search results, including:

- Find Cheap Airfares Fast**
www.kayak.com/Flights
kayak.com is rated ★★★★★
100s Of **Flights** On One Simple Site.
Compare Many Flying Options At Once
- Flight Amsterdam Dublin**
www.schipholtickets.nl/Ierland
Laatste stoelen met scherpe prijzen
Vergelijk & Bespaar snel online!
- Egypt Airlines**
www.egyptair.com/Egypt
Discover Luxury in the clouds
And enjoy our superior service
- Amsterdam Dublin?**
www.ebookers.nl/Dublin
Vliegticket Schiphol - **Dublin?**
Géén verborgen kosten. Boek nu!
- Van Amsterdam naar Dublin**
amsterdam-dublin.vlucht24.nl
Boek nu snel **Amsterdam - Dublin**.
Vliegtickets Vinden op Vlucht 24!
- Flights to Amsterdam**
www.vliegwinkel.nl/Amsterdam
All **flights & airlines** to **Amsterdam**
Compare, book and save money now!

Below the search results, there is a section titled "Flights from **Amsterdam, The Netherlands** (AMS) to **Dublin, Ireland**...". This section includes a blue airplane icon and the following information:

- Non-stop flights: 5 per day, 1h 40m duration
- Airlines: **Aer Lingus**
- + Schedule of non-stop flights

Below this section, there are several more search results, including:

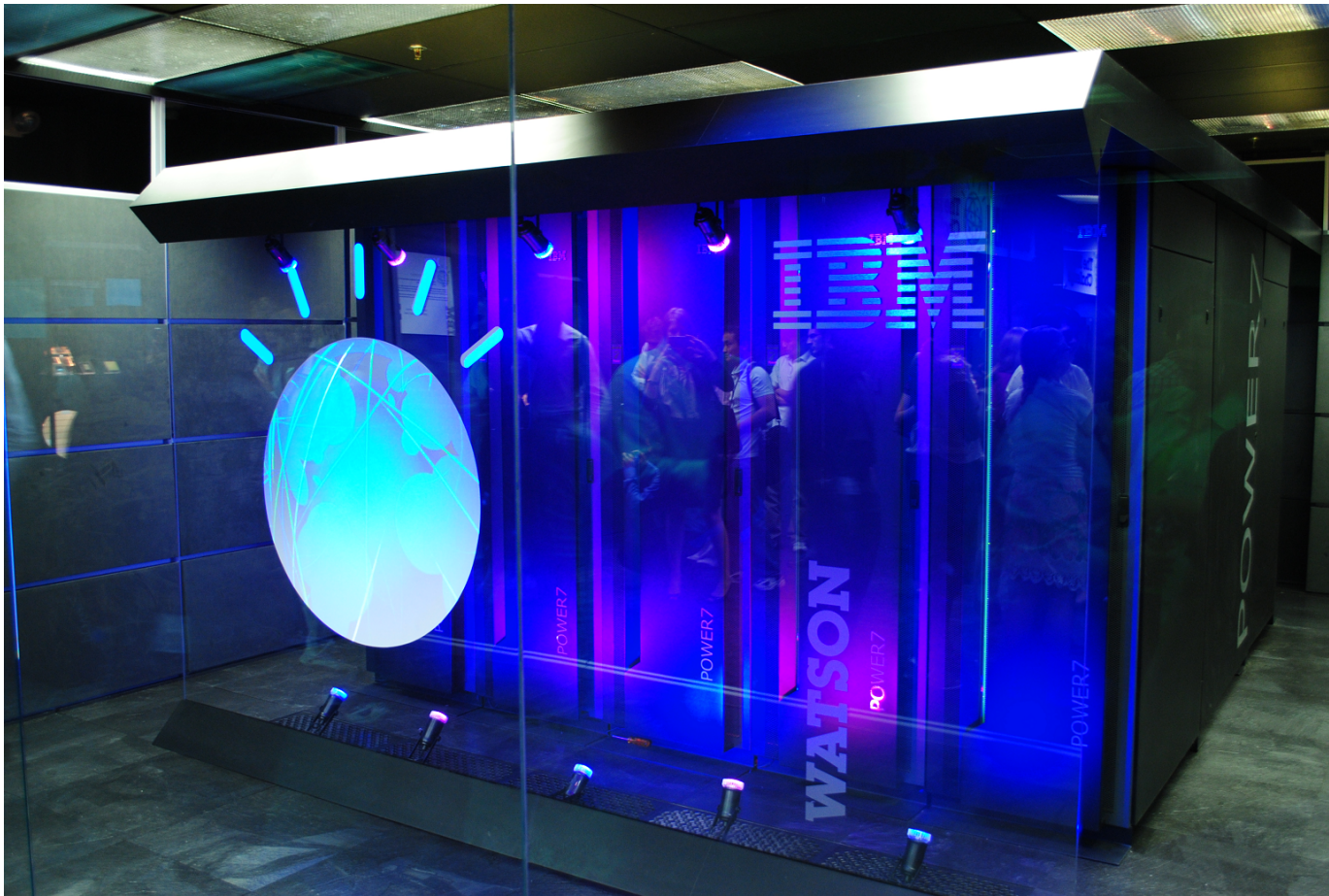
- Cheap Flights - To and from Dublin, Europe, UK & USA - Aer Lingus**
www.aerlingus.com/
Book cheap **flights** online today with Aer Lingus. Fly to Ireland, the UK, Europe and N. America including Canada with us as well as find hotels and more.
Book Flights - Web-Check-in - Manage Booking - Travel Information
- Cheap flights from Dublin to Amsterdam (Netherlands)**
www.cheapflights.co.uk/flights/Amsterdam/Dublin/
Dublin to Amsterdam flights. Search and compare cheap **flights** from **Dublin** (Ireland) to **Amsterdam** (Netherlands) to find the latest deals from all major airlines ...

On the left side of the search results, there is a blue callout box pointing to the search results area, containing the text "Query intent & vertical search".

On the right side of the search results, there is a blue callout box pointing to the search results area, containing the text "Ads adaptation".


Current systems and technologies III

Searching for answers (not documents): IBM Watson @ Jeopardy!




Current systems and technologies IV

Searching for expertise (not documents)



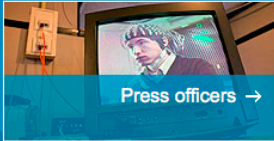
[Employees](#) | [Students](#) | [Contact](#) | [Working at](#) | [Nederlands / Dutch](#)


[EDUCATION](#) | [RESEARCH](#) | [NEWS AND EVENTS](#) | [ABOUT TILBURG UNIVERSITY](#) | [ALUMNI](#)

Google™ Aangepast zoeken 

Experts & Expertise

Many researchers, scientists and support staff are working at Tilburg University. Search for a certain individual, or within a certain expertise to find the right person.

[Press officers →](#)

[Mail the editorial staff →](#)

Experts

ABCDEFGHIJKLMNOPQRSTUVWXYZ

Search experts

Which researcher or scientist are you looking for?

☐ Also search for support staff

Expertise

ABCDEFGHIJKLMNOPQRSTUVWXYZ

Search expertise

In which field are you looking for a researcher or scientist?

Search results

Searched for: **information retrieval**

information retrieval

[M.M. van Zaanen](#)
[K. \(Kalliopi\) Zervanou](#)

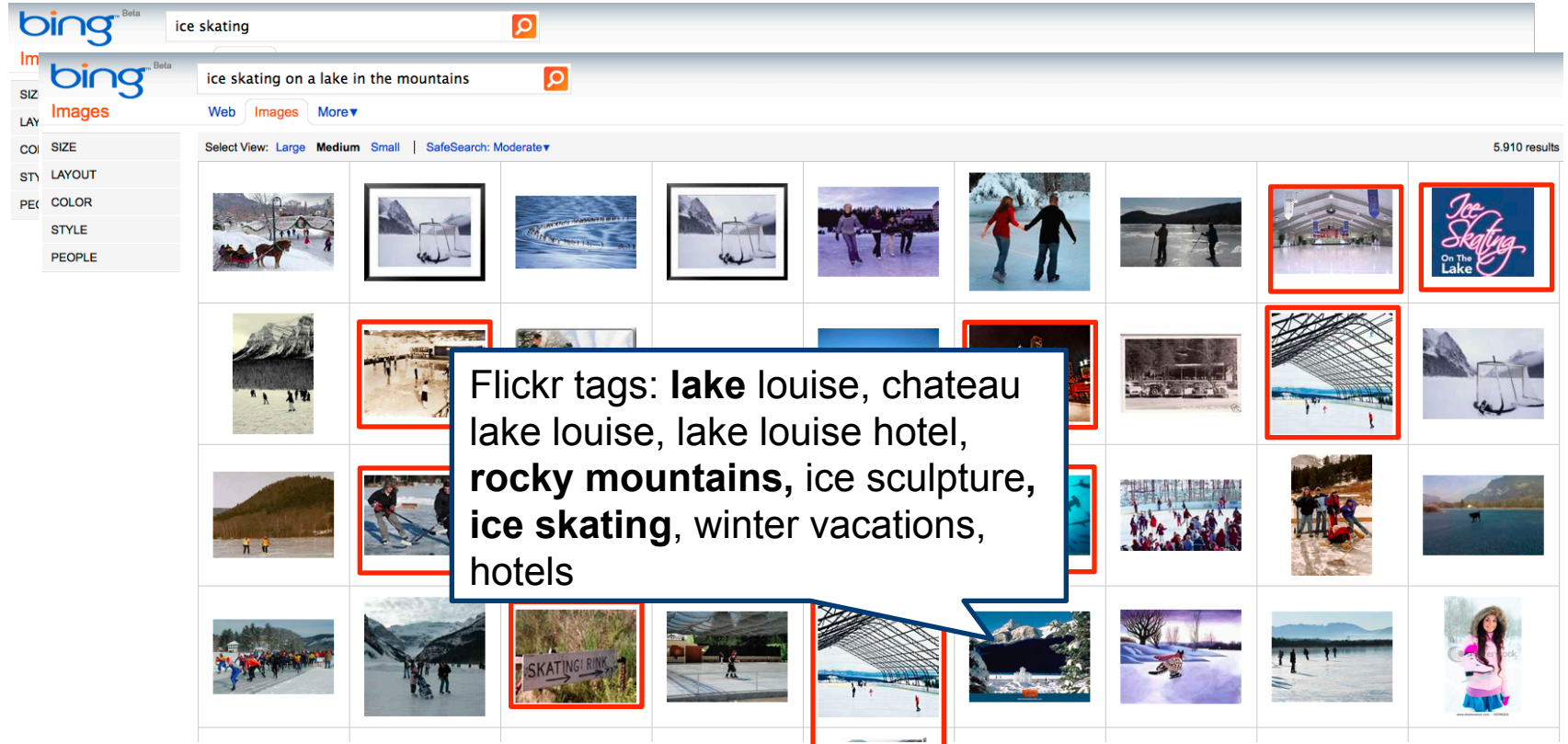
See also:

- [Digital library](#)
- [Documentary information](#)
- [Folksonomy](#)
- [Search engine](#)

[Search again](#)

<http://www.tilburguniversity.edu/webwijs/>

Current systems and technologies V




Current systems and technologies VI

Music retrieval

Music Ngram Viewer

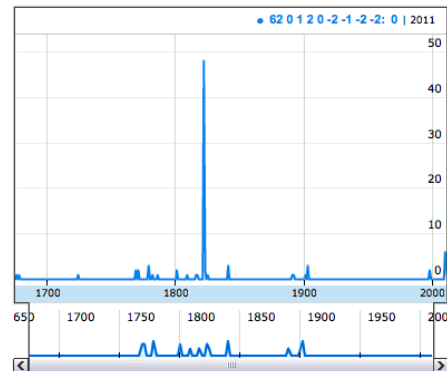
Please enter a melody or a sequence of **chords** (advanced use)



62 0 1 2 0 -2 -1 -2 -2 |

chord

Petrucchi Music Library Smoothing: 0 ☐ Normalized



query



results

filter search results (e.g. Mozart, winds, or quartet)

Symphony No.9
Beethoven, Ludwig van (1822)

You [YouTube](#) [score](#) pages
[18, 19,](#)
[25](#)

[score](#) pages [12, 14](#)
[12, 15,](#)
[28](#)

[score](#) pages [12, 14](#)

6 String Quartets, G.165-170 (Op.8)
Boccherini, Luigi (1769)

You [YouTube](#) [score](#) page [14](#)

String Quartets, Op.17
Haydn, Joseph (1771)

You [YouTube](#) [score](#) page [15](#)

Symphony No.33
Mozart, Wolfgang Amadeus (1779)

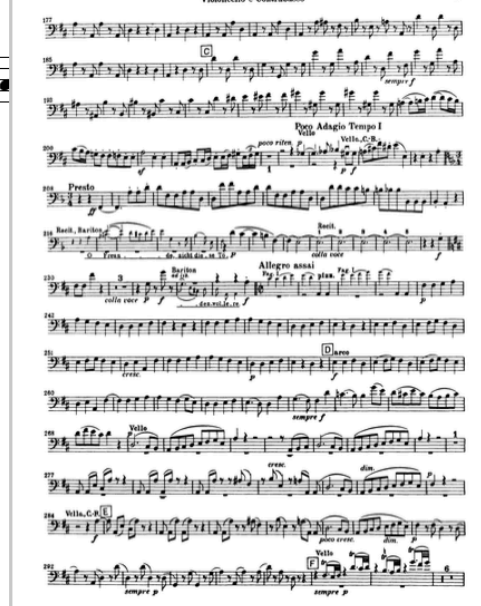
You [YouTube](#) [score](#) pages [12, 15](#)

Violin Sonata No.6
Beethoven, Ludwig van (1801)

You [YouTube](#) [score](#) page [4](#)

[next](#)
you can also browse using the chart

Bethoven — Symphony No. 9
Violoncello e Contrabbasso



Beethoven, Ludwig van: Symphony No.9
page 19

[previous page](#) [next page](#)

Run your own experiment! Raw data is available for download [here](#).

© 2011 Vladimir Viro - [About Music Ngram Viewer](#) - [Libraries](#) - [API](#) - [Contact](#) - [@Peachnote](#) on Twitter

<http://www.peachnote.com>

A few questions for you

- Who has programming experience in **Java**?
- Who is familiar with **functional programming**?
- Who has worked with **Hadoop**?
- Who has worked with **Amazon Web Services**?

Big data

What is 'Big data'?

- “Big data refers to enormous amounts of unstructured data produced by high-performance applications” [1]
 - Scientific computing applications
 - Social networks
 - E-government applications
 - Medical information systems
- Issues
 - Scalability
 - Heterogeneity
 - Data analysis

How big is big?

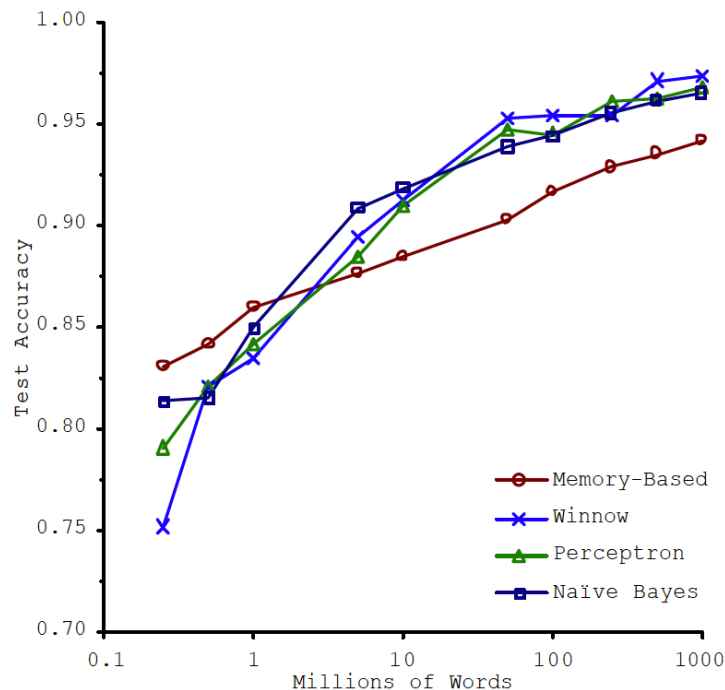
- Data repositories will only grow bigger with time.
- More data usually translates into more effective algorithms.

- YouTube: 4 billion views a day, one hour of video upload per second
- Facebook: 483 million daily active users (December 2011)
- Twitter: 140 million tweets on average per day (March 2011)
- Google: >1 billion searches per day (March 2011)
- Google processed 100 TB of data per day in 2004 and 20 PB data per day in 2008
- Large Hadron Collider at CERN: when fully functional, it will generate 15 petabytes of data per year
- Internet Archive: contains 2 petabytes of data, grows 20 terabytes per month (2011)

The more data, the better

The unreasonable effectiveness of data. A. Halevy, P. Norvig and F. Pereira, 2009.

“So, follow the data.”



Confusion set disambiguation

- *then* vs. *than*
- *to* vs. *two* vs. *too*
-

Scaling to very very large corpora for natural language disambiguation. M. Banko and E. Brill, 2001.

Cloud computing

- “Anything running inside a browser that gathers and stores user-generated content” (Jimmy Lin)

- Utility computing

- A computing resource
- A “cloud user” buys from a “cloud provider” (pay-as-you-go)
 - Virtual machine instances
- IaaS: infrastructure as a service
- Amazon Web Services is the dominant provider

Region: EU (Ireland)	Linux/UNIX Usage	Windows Usage
Standard On-Demand Instances		
Small (Default)	\$0.095 per hour	\$0.12 per hour
Large	\$0.38 per hour	\$0.48 per hour
Extra Large	\$0.76 per hour	\$0.96 per hour
Micro On-Demand Instances		
Micro	\$0.025 per hour	\$0.035 per hour
Hi-Memory On-Demand Instances		
Extra Large	\$0.57 per hour	\$0.62 per hour
Double Extra Large	\$1.14 per hour	\$1.24 per hour
Quadruple Extra Large	\$2.28 per hour	\$2.48 per hour
Hi-CPU On-Demand Instances		
Medium	\$0.19 per hour	\$0.29 per hour
Extra Large	\$0.76 per hour	\$1.16 per hour

Amazon EC2 Pricing, 6th February 2012

Cloud computing

- “Anything running inside a browser that gathers and stores user-generated content” (Jimmy Lin)
- Utility computing
 - A computing resource as a metered service
 - A “cloud user” buys any amount of computing power from a “cloud provider” (pay-per-use)
 - Virtual machine instances
 - IaaS: infrastructure as a service
 - Amazon Web Services (EC2: elastic compute cloud) is the dominant provider

MapReduce

- ① Programming model for distributed computations on **large-scale** data, inspired by the functional programming paradigm
- ② Execution framework for clusters of commodity hardware
- Developed by researchers at Google in 2003
 - Built on principles in parallel and distributed processing
- “MapReduce is used for the generation of data for Google’s production web search service, for sorting, for data mining, for machine learning and many other systems” [12]
- Designed for **batch** processing over large data sets

Ideas behind MapReduce I

- Scale “out”, not “up”
 - Many commodity servers are more cost effective than few high-end servers
- Assume failures are common
 - A 10,000-server cluster with a mean-time between failures of 1000 days experiences on average 10 failures a day.
- Move processes to the data
 - Moving the data around is expensive
- Process data sequentially and avoid random access
 - Data sets do not fit in memory, disk-based access (slow)

Ideas behind MapReduce II

- Hide system-level details from the application developer
 - Frees the developer to think about the task at hand only (no need to worry about deadlocks, ...)
 - MapReduce takes care of the system-level details (separation of what and how to compute)

To be continued in the next lecture

Sources

- <http://www.flickr.com/photos/annarbor/4349874305>
- <http://www.flickr.com/photos/annarbor/4349874915>
- <http://www.flickr.com/photos/ccacnorthlib/4775120660>
- <http://www.flickr.com/photos/awhitis/4266617452/>
- <http://www.flickr.com/photos/karenhorton/4893479758>
- http://en.wikipedia.org/wiki/File:IBM_Watson.PNG



- ① On the history of evaluation in IR. Stephen Robertson. 2008
- ② Information retrieval and digital libraries: lessons of research. Karen Spaerck Jones. 2006.
- ③ Information retrieval viewed as temporal signaling. Calvin Mooers. 1950.
- ④ Making information retrieval pay. American Chemical Society 118th Meeting, abstracts of papers
- ⑤ Information retrieval. Keith van Rijsbergen
- ⑥ Experimentation as a way of life: Okapi at TREC. S. Robertson, S. Walker and M. Beaulieu, 2000.
- ⑦ <http://www.cs.cornell.edu/Info/Department/Annual96/Beginning/salton.html>
- ⑧ As we may think. Vannevar Bush, Atlantic Monthly, 1945. <http://www.theatlantic.com/magazine/archive/1945/07/as-we-may-think/3881/>
- ⑨ Google Effects on Memory: Cognitive Consequences of Having Information at Our Fingertips. B. Sparrow, J. Liu and M. Wegner, Science, 2011. <http://www.sciencemag.org/content/early/2011/07/13/science.1207745>