

Information Retrieval

Dhruv Gupta

dhruv.gupta@ntnu.no

23-August-2022



NTNU

|

Norwegian University of
Science and Technology

1 Course Information

- Course Team
- Course Prerequisites
- Course References
- Guest Lectures
- Assignments
- Examination
- Course Contents
- Lecture Plan

2 Introduction

- Definition
- Search in Libraries
- Search in Digital Libraries
- Searching the Web

Course Team — Instructors

- I will cover the **first half of the course**.
- Heri will cover the **second half of the course** (partly via recorded lectures).



Dhruv Gupta (*dhruv.gupta@ntnu.no*)



Heri Ramampiaro (*heri@ntnu.no*)

Course Team — Teaching Assistants

- **Teaching Assistant Leader:** David Baumgartner and Shiva Shadrooh.
- **Teaching Assistants:**
 - 1 Anthony Vu (*thvu@stud.ntnu.no*)
 - 2 Agnes Hjeltnes (*agnesgh@stud.ntnu.no*)



David Baumgartner
(*david.baumgartner@ntnu.no*)



Shiva Shadrooh
(*shiva.shadrooh@ntnu.no*)

Course Team — General

- **Course ID:** TDT4117 Information Retrieval
- Email to Teaching Assistants: *tdt4117-assist@idi.ntnu.no.*
- **Lectures:** Tuesdays (10:15-12:00 HRS).
- **Location:** VE1 Verkstedteknisk.
- **Piazza** is the main platform for discussions.
Join here: *piazza.com/ntnu.no/fall2022/tdt4117.*
- **Blackboard** is the main platform for course-related information.

Course Prerequisites

- The subject material is at an **advanced level**.
- Recommended prerequisites:
 - IT1104 Programming Advanced Course
 - TDT4100 Object-Oriented Programming
 - TDT4145 Data Modelling, Databases and Database Management Systems

Course References

- Main course textbook:
Baeza-Yates and Ribeiro-Neto,
“Modern Information Retrieval”,
Second Edition.
Pearson Education Limited, 2011.

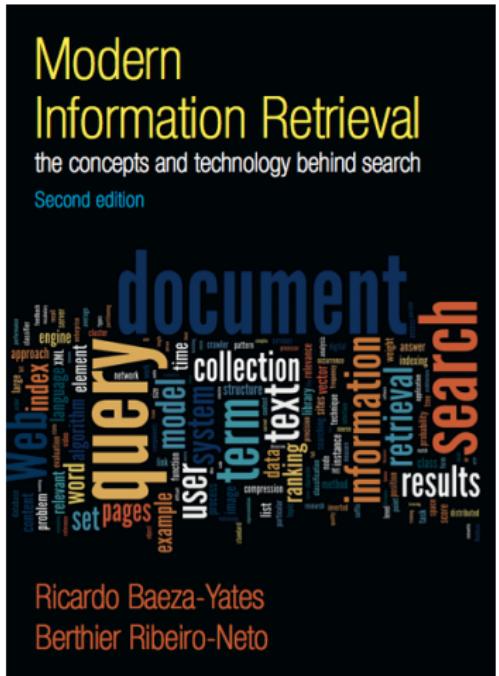


Image Credit: <http://grupoweb.upf.es/mir2ed/>

Course References

- Reference book: Manning et al., “Introduction to Information Retrieval”, First Edition. Cambridge University Press, 2008.
- Book is publicly available at: [https://nlp.stanford.edu/IR-book/.](https://nlp.stanford.edu/IR-book/)

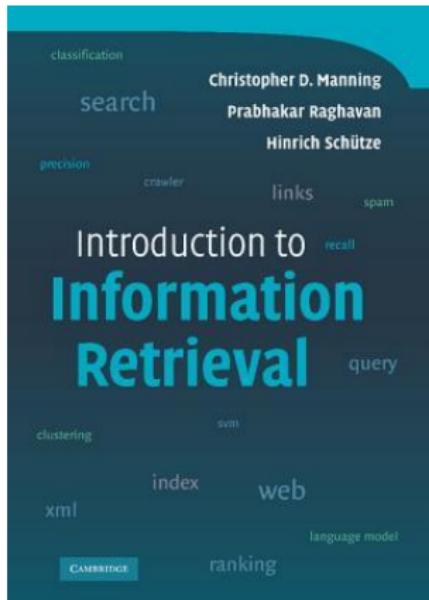


Image Credit: <https://www.goodreads.com/book/show/3278309-introduction-to-information-retrieval>

Guest Lectures

- Tentatively, there will be a guest lecture from a speaker in the Industry.
- Intention is to supplement the course knowledge from industry by daily practitioners.
- Announcement at an upcoming date.

Assignments

- 5 mandatory exercises, where 4 out of 5 must be approved.
- Submissions can be made either individually or by two-person groups.
- The exercises will be published approximately two weeks before the submission deadline.
- We expect that most of the exercises are done correctly, otherwise one an improved submission will be required.
- Plagiarism will lead to fail grade on the assignment.
- Tentative exercise plan:

Available (week)	Deadline (week)	Theme
08.09.2022 (36)	22.09.2022 (38)	Introductory Topics and Similarity Models
22.09.2022 (38)	06.10.2022 (40)	Language Models and Evaluation in IR
06.10.2022 (40)	20.10.2022 (42)	Text Operations
20.10.2022 (42)	03.11.2022 (44)	Indexing
03.11.2022 (44)	17.11.2022 (46)	Web Search

Examination

- Examination date: 07.December.2022 at 15:00.
- It will be a **written exam with grades** accounting for 100% of the grade.
- **Old examinations:** posted on Blackboard.
There is also a solution sketch for most questions.

Course Contents¹

● Academic Content

- The course concerns **automatic document storage and retrieval**. In this case, the term document includes sounds and images as well as text.

With this course you will learn about:

- 1 File Organization,
- 2 Query Operations,
- 3 Document Operations,
- 4 Knowledge-based Textual Information Retrieval, and
- 5 Multimedia Information Retrieval.

● Learning Objectives

- The students will learn and understand the principle, techniques and methods behind information retrieval.

¹From: <https://www.ntnu.edu/studies/courses/TDT4117/2022>

Lecture Plan (Tentative)

Week/ Day	Topic	Lecturer	Chapter Number
34 / Tuesday	Welcome, General Information, and Introduction	Dhruv Gupta	1
35 / Tuesday	Classical Similarity Models	Dhruv Gupta	3
36 / Tuesday	Classical Similarity Models (continued), BM25, and Language Model	Dhruv Gupta	3
37 / Tuesday	Classical Similarity Models (continued), BM25, and Language Model	Dhruv Gupta	3
38 / Tuesday	Evaluation in Information Retrieval	Dhruv Gupta	4
39 / Tuesday	User Relevance Feedback and Query Expansion	Dhruv Gupta	5
-	Text Operation	Heri Ramampiaro (recorded video)	6
-	Indexing and Searching	Heri Ramampiaro (recorded video)	9
-	Guest Lectures (Tentative)	Tentative	-
-	Web Search and Search Engines	Heri Ramampiaro (recorded video)	11
-	Introduction to Multimedia Retrieval and Image Retrieval	Heri Ramampiaro (recorded video)	-
-	Audio Retrieval	Heri Ramampiaro (recorded video)	-
-	Audio Retrieval	Heri Ramampiaro (recorded video)	-

1 Course Information

- Course Team
- Course Prerequisites
- Course References
- Guest Lectures
- Assignments
- Examination
- Course Contents
- Lecture Plan

2 Introduction

- Definition
- Search in Libraries
- Search in Digital Libraries
- Searching the Web

1 Course Information

- Course Team
- Course Prerequisites
- Course References
- Guest Lectures
- Assignments
- Examination
- Course Contents
- Lecture Plan

2 Introduction

● Definition

- Search in Libraries
- Search in Digital Libraries
- Searching the Web

Information Retrieval — Definition

“Information retrieval (IR) is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers).”²

²Manning et al., “Introduction to Information “Retrieval”, First Edition. Cambridge University Press, 2008.

1 Course Information

- Course Team
- Course Prerequisites
- Course References
- Guest Lectures
- Assignments
- Examination
- Course Contents
- Lecture Plan

2 Introduction

- Definition
- **Search in Libraries**
- Search in Digital Libraries
- Searching the Web

Search in Libraries

- Information Retrieval \equiv Search.
- Search before digitization: libraries.



https://commons.wikimedia.org/wiki/File:Bookshelf_at_Yale.jpg

Search in Libraries

- Search before digitization: libraries.
- Search procedure using catalogs.
- Catalogs contained manually curated metadata about the book.



Sample Catalog Record
Author: Kesey, Ken.
Title: One flew over the cuckoo's nest, a novel.
Published: New York, Viking Press [1962]
LC Call No.: PZ4.K42On
Subjects: Psychiatric hospital patients--United States--fiction
Control No.: 62008602

https://commons.wikimedia.org/wiki/File:Sample_Catalog_Record.png

https://en.wikipedia.org/wiki/Library_catalog#/media/File:Schlagwortkatalog.jpg

Search in Libraries

- Search before digitization: libraries.
- Search procedure using manually created catalogs.
- **Location** of book encoded using **Dewey Decimal Classification (DDC)**.

Decoding Dewey Decimal Call Numbers		
Main Class	800	Literature
Division	810	American literature in English
Section	813	American fiction in English
	813.54	... further narrowing of topic
	813.54 M37	Cutter Number identifying author's name
	813.54 M37 2007	Edition date

Digitizing Libraries

- Converting physical objects of knowledge to digital counterparts.
- Famous initiatives: Project Gutenberg (<https://gutenberg.org/>) and Internet Archive (<https://archive.org/>).

The image shows two screenshots of library digitization websites side-by-side.

Project Gutenberg (Left):

- Header:** "Project Gutenberg" logo, "About", "Search and Browse", "Help", "Donation", "Feedback".
- Section:** "Welcome to Project Gutenberg".
- Text:** "Project Gutenberg is a library of over 60,000 free eBooks".
- Text:** "Choose among free epub and Kindle eBooks, download them or read them online. You will find the works for which U.S. copyright has expired. Thousands of volumes digitized and diligently proofread by volunteers".
- Image:** A grid of book covers including "Afghan Book No. 289: Afghan", "A Son of Ishmael", "Young Engineer's Guide", "A Gloucestershire Lad at Home", "Kuulemani mmevottaa", and "Tales of the Wild and the Wonderful".
- Text:** "Some of our latest eBooks Click Here for more latest books".
- Text:** "50 years of eBooks 1971-2021. The first eBook for reading enjoyment and unlimited free redistribution is here. Read more about this listing (involves). Project Gutenberg is grateful to all volunteers who Gutenberg offers a vibrant and growing collection of the world's great literature. Read, enjoy, and savor".
- Text:** "No fee or registration! Everything from Project Gutenberg is gratis, fore, and completely without fees. Please consider a small donation to help Project Gutenberg digitize more books, maintain its online offerings. Other ways to help include digitizing, proofreading and formatting, or reporting errors".
- Text:** "No special apps needed! Project Gutenberg eBooks require no special apps to read, just the regular computer and reading device. There have been reports of sites that charge fees for customizations from Project Gutenberg. Some of the apps might have worthwhile features, but none are required".
- Section:** "Find Free eBooks".
- List:**
 - [Search and browse](#): By author, title, subject, language, type, popularity, and more.
 - [Bookshelves](#): of related eBooks.
 - [Frequently downloaded](#): Top 100, or ranked by popularity.
 - [Offline catalogs](#): handy eBook listings and metadata to consult offline.
 - [Recently added](#): The latest new and updated eBooks.
- May only access these e-books for non-commercial purposes.

Internet Archive (Right):

- Header:** "INTERNET ARCHIVE" logo, "25", "ABOUT", "BLOG", "PROJECTS", "HELP", "DONATE", "CONTACT", "JOBS", "VOLUNTEER", "PEOPLE".
- Text:** "Search the history of over 600 billion web pages on the Internet".
- Section:** "WayBack Machine".
- Text:** "Internet Archive is a non-profit library of millions of free books, movies, software, music, websites, and more".
- Image:** A blue banner with "25 Years From Wayback to Way forward" and a link to "Celebrating 25".
- Text:** "Search" and "Advanced Search" fields.
- Section:** "Announcements".
- Text:** "Performance as the Internet Archive turns 25", "Internet Archive Launches Collaborative, Web-Based Art History Preservation and Access Initiative", "Drexler Kudos named to the Library of Congress' Copyright Public Motorization Committee".
- Section:** "Top Collections at the Archive".
- Grid:** A 4x5 grid of collection icons and names:
 - American Libraries (9,076 items)
 - LibriVox (15,487 items)
 - Canadian Libraries (490,400 items)
 - Live Music Archive (297,759 items)
 - Electric Sheep (453 items)
 - APK Archive (76,345 items)
 - Folksonomy: A Library of... (802,152 items)
 - University of Toronto (216,962 items)
 - The Phone Software... (102,799 items)
 - California Digital Library (160,833 items)
- Terms of Service last updated 12/31/2016

1 Course Information

- Course Team
- Course Prerequisites
- Course References
- Guest Lectures
- Assignments
- Examination
- Course Contents
- Lecture Plan

2 Introduction

- Definition
- Search in Libraries
- Search in Digital Libraries**
- Searching the Web

Search in Digital Libraries

- Since, most catalogs are in **structured format** we can leverage a **database solution**.
- Naïve approach:** map existing manual catalogs to **digital database schema**.

ID	Author	Book Name	ISBN	Subject
10014	Ricardo Baeza-Yates	Modern Information Retrieval	9780321416919	Computer Science

Search in Digital Libraries

- Since, most catalogs are in **structured format** we can leverage a **database solution**.
- Naïve approach:** map existing manual catalogs to **digital database schema**.

ID	Author	Book Name	ISBN	Subject
10014	Ricardo Baeza-Yates	Modern Information Retrieval	9780321416919	Computer Science

```
SELECT *
FROM  DIGITAL_LIBRARY
WHERE BOOK_NAME LIKE '%Information Retrieval%'
```

Search in Digital Libraries

- Since, most catalogs are in **structured format** we can leverage a **database solution**.
- **Naïve approach:** map existing manual catalogs to **digital database schema**.
- Drawbacks:
 - Slow
 - Unstructured text within documents is in-accessible
 - **Presentation of results**
 - Determination of which is the best book

Search in Digital Libraries

- Since, most catalogs are in **structured format** we can leverage a **database solution**.
- **Naïve approach:** map existing manual catalogs to **digital database schema**.
- Drawbacks:
 - Slow
 - Unstructured text within documents is in-accessible
 - Presentation of results
 - Determination of which is the best book
- Can not cast "**Information Retrieval**" as "**Data Retrieval**" task!

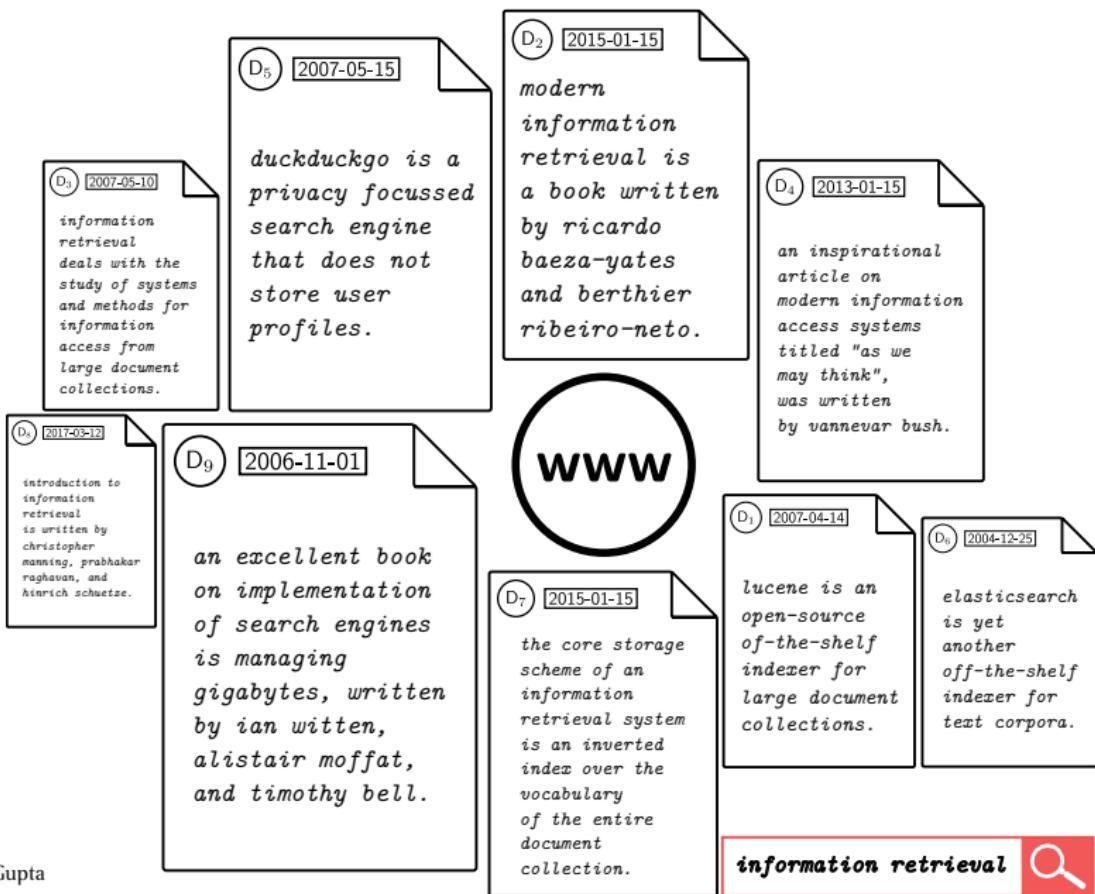
1 Course Information

- Course Team
- Course Prerequisites
- Course References
- Guest Lectures
- Assignments
- Examination
- Course Contents
- Lecture Plan

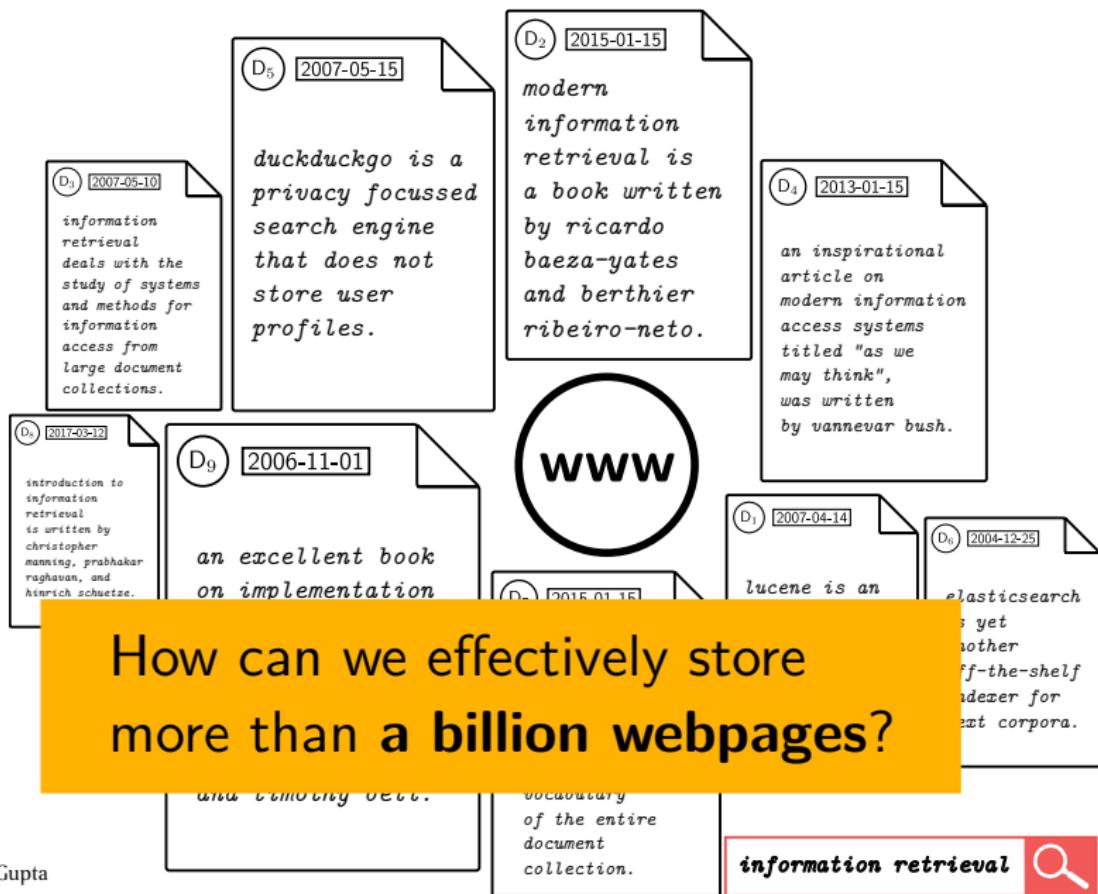
2 Introduction

- Definition
- Search in Libraries
- Search in Digital Libraries
- Searching the Web**

Searching the Web — Document Collection



Searching the Web — Challenge of Size



Searching the Web — Fault Tolerance and Latency



<https://duckduckgo.com/traffic>

Searching the Web — Users, Queries, and Interface

information retrieval



Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal for knowledge about various topics.
Wikipedia entry of the topic web search.

Modern Information Retrieval

<http://grupoweb.upf.es/mir2ed/>

Modern information retrieval is a book written by Ricardo
Baeza-Yates and Berthier Ribeiro-Neto.

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning,
Prabhakar Raghavan and Hinrich Schütze,

Information Retrieval

" Information retrieval is the process of obtaining information system resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other content-based indexing."

- Wikipedia

https://en.wikipedia.org/wiki/Information_retrieval

Searching the Web — Queries and Information Needs

information retrieval



Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal fo

Wikipedia entry of the topic w

Modern Information Re

<http://grupoweb.upf.es/mir2/>

Modern information retrieval is a book written by Ricardo

Baeza-Yates and Berthier Ribeiro-Neto.

Information Retrieval

" Information retrieval is the process of obtaining information system resources that are relevant to an information need from a collection

Users express their **information needs**
using natural language expressions.

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning,

Prabhakar Raghavan and Hinrich Schütze,

Searching the Web — Queries

information retrieval



Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal fo

Wikipedia entry of the topic w

Modern Information Re

<http://grupoweb.upf.es/mir2/>

Modern information retrieval is a book written by Ricardo

Baeza-Yates and Berthier Ribeiro-Neto.

Information need can be expressed as:
keywords, phrases, and many more!

Information Retrieval

" Information retrieval is the process
of obtaining information system
resources that are relevant to an
information need from a collection

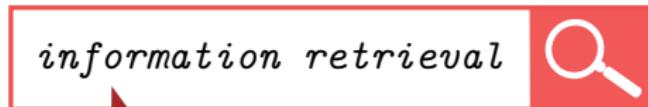
Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning,

Prabhakar Raghavan and Hinrich Schütze,

Searching the Web — Query Operators



Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal fo

Wikipedia entry of the topic w

Modern Information Re

<http://grupoweb.upf.es/mir2/>

Modern information retrieval is a book written by Ricardo

Baeza-Yates and Berthier Ribeiro-Neto.

Information Retrieval

" Information retrieval is the process of obtaining information system resources that are relevant to an information need from a collection

**Boolean Operators further enable
conjunctive, disjunctive and
negation queries!**

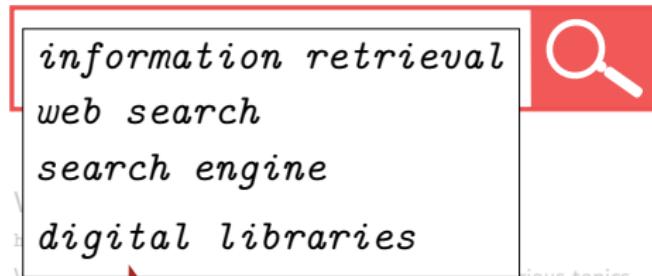
Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning,

Prabhakar Raghavan and Hinrich Schütze,

Searching the Web — Query Reformulations



Wikipedia entry of the topic web search.

Wikipedia entry of the topic web search.

Modern Information Retrieval

<http://grupoweb.upc.es/mir2ed/>

Modern information retrieval is a book by

Baeza-Yates and Berthier Ribeiro-Neto.

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/index.html>

Introduction to Information Retrieval is written by Christopher Manning,

Prabhakar Raghavan and Hinrich Schütze,

Information Retrieval

" Information retrieval is the process of obtaining information system resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other content-based indexing."

- Wikipedia

https://en.wikipedia.org/wiki/Information_retrieval

**Query reformulations can
be identified using
pseudo-relevance feedback!**

Searching the Web — Full-Text Search

information retrieval



Search results are identified by **full-text search** through document contents.

Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal for knowledge about various topics.

Wikipedia entry of the topic **web search**.

of obtaining information system resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other content-based indexing."

- Wikipedia

https://en.wikipedia.org/wiki/Information_retrieval

Modern Information Retrieval

<http://grupoweb.upf.es/mir2ed/>

Modern **information retrieval** is a book written by Ricardo Baeza-Yates and Berthier Ribeiro-Neto.

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to **Information Retrieval** is written by Christopher Manning, Prabhakar Raghavan and Hinrich Schütze,

Searching the Web — Synonym Discovery

information retrieval



Search results are also identified by **synonym discovery!**

Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal for knowledge about various topics.

Wikipedia entry of the topic **web search**.

of occurring information system resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other content-based indexing.”

Modern Information Retrieval

<http://grupoweb.upf.es/mir2ed/>

Modern information retrieval is a book written by Ricardo

Baeza-Yates and Berthier Ribeiro-Neto.

- Wikipedia

https://en.wikipedia.org/wiki/Information_retrieval

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning,

Prabhakar Raghavan and Hinrich Schütze,

Searching the Web — Relevance

information retrieval



Search results are
ranked by relevance
to the end-user!

Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal for knowledge about various topics.

Wikipedia entry of the topic web search.

Modern Information Retrieval

<http://grupoweb.upf.es/mir2ed/>

Modern information retrieval is a book written by Ricardo Baeza-Yates and Berthier Ribeiro-Neto.

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning, Prabhakar Raghavan and Hinrich Schütze,

resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other content-based indexing."

- Wikipedia
https://en.wikipedia.org/wiki/Information_retrieval

Searching the Web — Personalization

information retrieval



The relevance can further be "improved" based on **previous searches by the user!**

Modern Information Retrieval

<http://grupoweb.upf.es/mir2ed/>

Modern information retrieval is a book written by Ricardo Baeza-Yates and Berthier Ribeiro-Neto.

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning, Prabhakar Raghavan and Hinrich Schütze,

resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other concepts indexing."

– Wikipedia
[wiki/Information_retrieval](#)

Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal for knowledge about various topics.

Wikipedia entry of the topic web search.

Searching the Web — Machine Learning

information retrieval



The relevance can further be "improved" based on previous searches by the user and **machine learning!**

Modern Information Retrieval

<http://grupoweb.upf.es/mir2ed/>

Modern information retrieval is a book written by Ricardo Baeza-Yates and Berthier Ribeiro-Neto.

resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other components indexing."

$$f(x)$$

- Wikipedia
https://en.wikipedia.org/w/index.php?title=Information_retrieval&oldid=98301110

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning, Prabhakar Raghavan and Hinrich Schütze,

Web Search - Wikipedia

https://en.wikipedia.org/wiki/Search_engine

Wikipedia is a public portal for knowledge about various topics.

Wikipedia entry of the topic web search.

Searching the Web — User Interfaces

information retrieval



User Interface:

To address the increasingly complicated information needs of users, the search interface is constantly evolving.

Modern Information Retrieval

<http://grupoweb.upf.es/mir2ed/>

Modern information retrieval is a book written by Ricardo Baeza-Yates and Berthier Ribeiro-Neto.

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning, Prabhakar Raghavan and Hinrich Schütze,

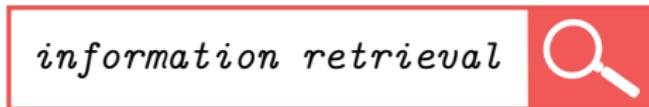
Information Retrieval

" Information retrieval is the process of obtaining information system resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other content-based indexing."

- Wikipedia

https://en.wikipedia.org/wiki/Information_retrieval

Searching the Web — User Interfaces



Modern Information Retrieval

<http://grupoweb.upf.es/mir2ed/>

Modern information retrieval is a book written by Ricardo Baeza-Yates and Berthier Ribeiro-Neto.

Introduction to Information Retrieval

<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

Introduction to Information Retrieval is written by Christopher Manning, Prabhakar Raghavan and Hinrich Schütze,

Information Retrieval

" Information retrieval is the process of obtaining information system resources that are relevant to an information need from a collection of those resources. Searches can be based on full-text or other content-based indexing."

- Wikipedia

https://en.wikipedia.org/wiki/Information_retrieval

Shopping User Interface

Alternative user interface allows users to buy related items to their queries instantly !

Dodge Charger
<http://grupoweb.upf.es/mir2ed/>
<https://www.gutenberg.org/cache/epub/123/pg123.html>
<https://www.gutenberg.org/cache/epub/123/pg123.html>

Searching the Web — Summary

- **Information retrieval:** making knowledge accessible to the user.
- **Objectives:**
 - Fast (millisecond response time) and fault tolerant.
 - Full-text search.
 - Ordering of documents according to relevance to the user.
 - User-centric presentation of results.