# Modern Information Retrieval

## Introduction

Hamid Beigy

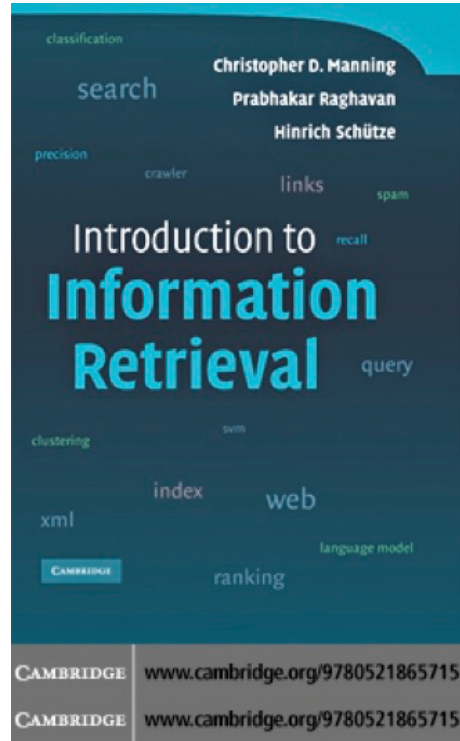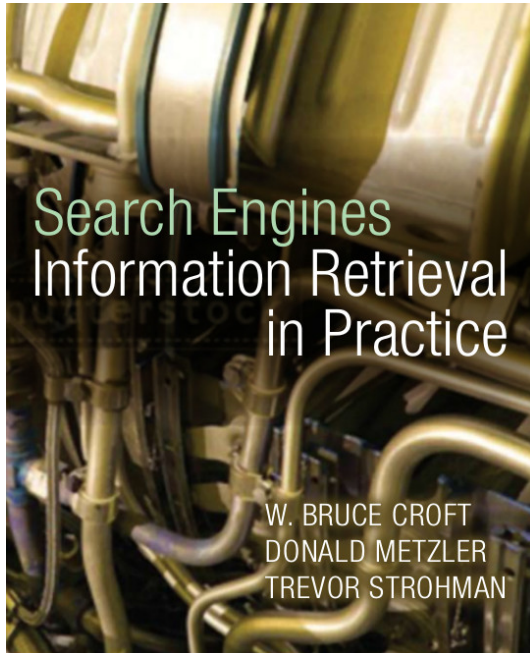Sharif University of Technology

1403-11-20

# Course Information

1. Course name : **Modern Information Retrieval**

2. Instructor : Hamid Beigy     Email : beigy@sharif.edu

3. Class : CE 201

4. Virtual class link: https://vc.sharif.edu/beigy

5. Course Website: http://sharif.edu/~beigy/14032-40324.html

6. Lectures: Sat-Mon (9:00-10:30)

7. Teaching Assistant : Reza Tavakoli     Email: seyedreza.shiyade@gmail.com

- Evaluation:

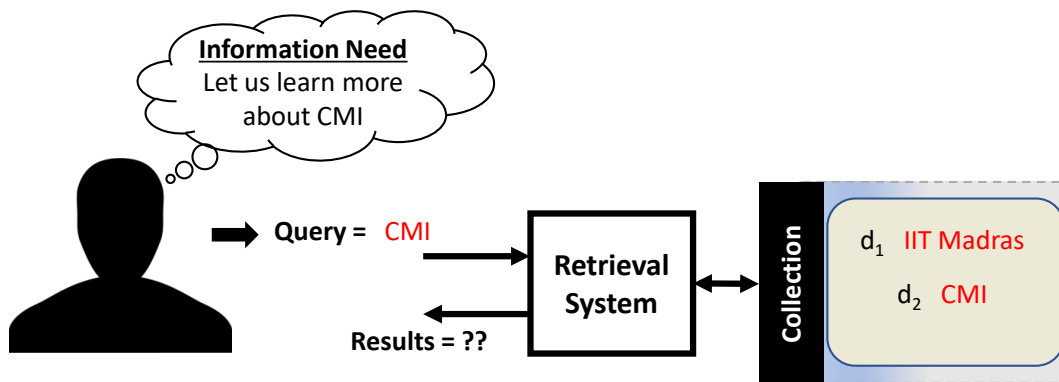| | | |
|---|---|---|
| Mid-term exam | 25% | 1404-01-30 |
| Final exam | 30% | |
| Practical Assignments | 35% | |
| Quiz | 15% | |

## References

Baeza-Yates, Ricardo and Berthier Ribeiro-Neto (2011). *Modern Information Retrieval*. 2nd. USA: Addison-Wesley Publishing Company. ISBN: 9780321416919.

Croft, W. Bruce, Donald Metzler, and Trevor Strohman (2009). *Search Engines - Information Retrieval in Practice*. Pearson Education.

Kowalski, Gerald (2010). *Information Retrieval Architecture and Algorithms*. 1st. Berlin, Heidelberg: Springer-Verlag. ISBN: 1441977155, 9781441977151.

Li, Hang (2011). *Learning to Rank for Information Retrieval and Natural Language Processing*. Morgan & Claypool Publishers.

Manning, Christopher D., Prabhakar Raghavan, and Hinrich Schütze (2008). *Introduction to Information Retrieval*. New York, NY, USA: Cambridge University Press.

Mitra, Bhaskar and Nick Craswell (2018). "An Introduction to Neural Information Retrieval". In: *Foundations and Trends in Information Retrieval* 13.1, pp. 1–126.

# Introduction

**Life without search engines is difficult to imagine!**

**Definition (Information retrieval )**

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).

1. If you know
   - Which stock to invest in?
   - Which faculty to work with?
   - How to get into a top college?
   - Which course to register for?
   - What to study?
   - How to prepare for job interviews?
2. **If only you had the information, you could rule this world.**
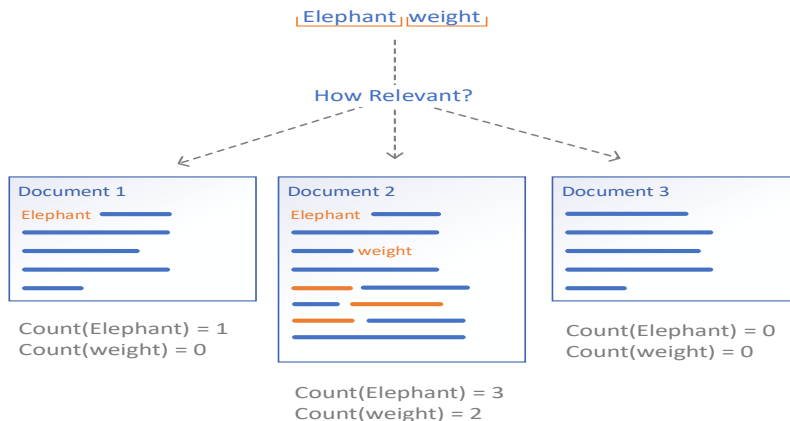3. **What happens when you lose access to all your information?**

1. Document Collection: units we have built an IR system over such as

   - memos
   - book chapters paragraphs
   - scenes of a movie
   - turns in a conversation...

2. Some applications

   - E-mail search
   - Searching your laptop
   - Corporate knowledge bases
   - Legal information retrieval

3. An information need is the topic about which the user desires to know more about.

4. A query is what the user conveys to the computer in an attempt to communicate the information need.

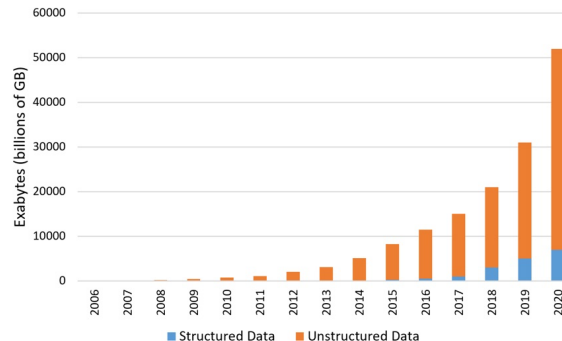1. A document is relevant if the user perceives that it contains information of value with respect to their personal information need.

2. Are the retrieved documents
   - about the target subject ?
   - up-to-date?
   - from a trusted source?
   - satisfying the user's needs?

3. How should we rank documents in terms of these factors?

1. The effectiveness of an IR system is determined by two key statistics:

   - Precision: What fraction of returned results are relevant to information need?

   - Recall: What fraction of relevant documents in collection were returned by system?

   - What is the best balance between the two?

     - Easy to get perfect recall: just retrieve everything

     - Easy to get good precision: retrieve only the most relevant

Elephant weight

How Relevant?

Document 1
Elephant
Count(Elephant) = 1
Count(weight) = 0

Document 2
Elephant
weight
Count(Elephant) = 3
Count(weight) = 2

Document 3
Count(Elephant) = 0
Count(weight) = 0

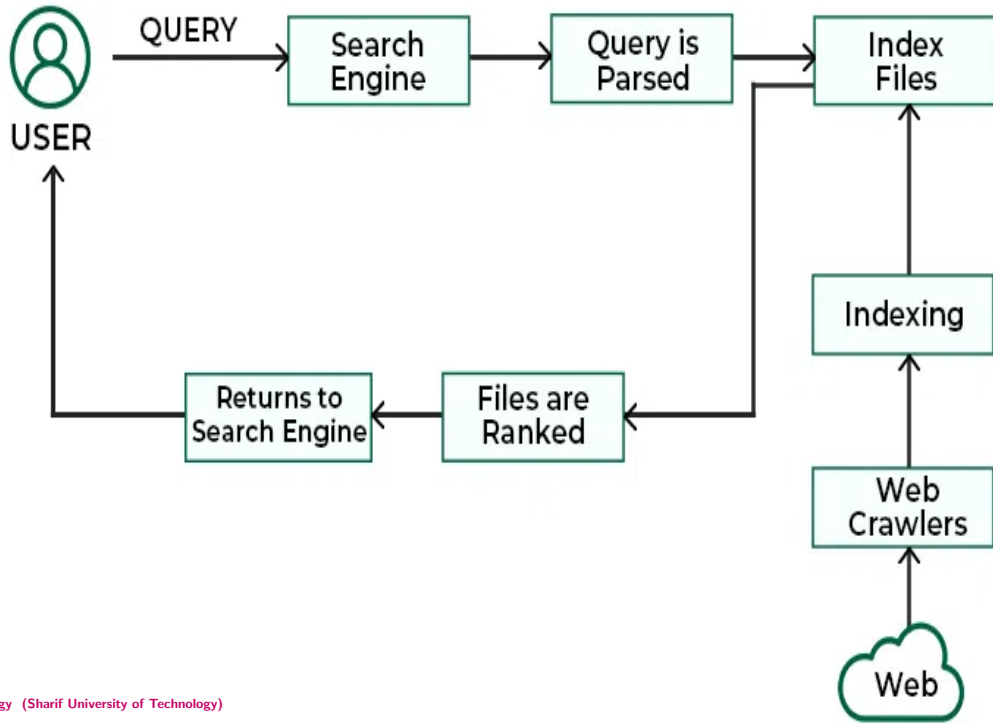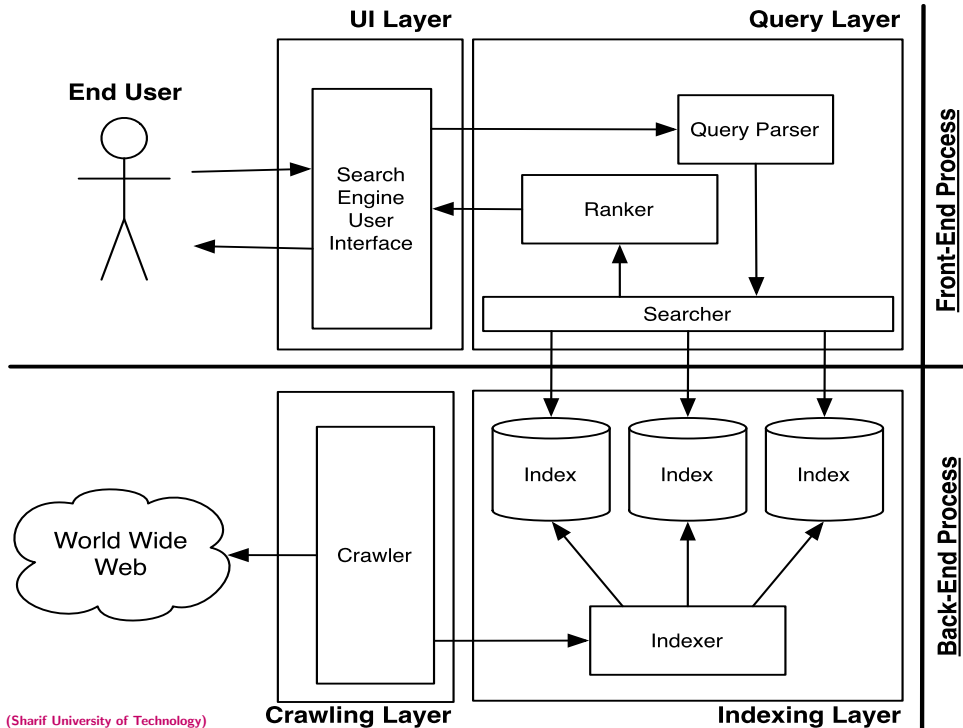1. Unstructured data: a formal, semantically overt, easy-for-computer structure is missing.
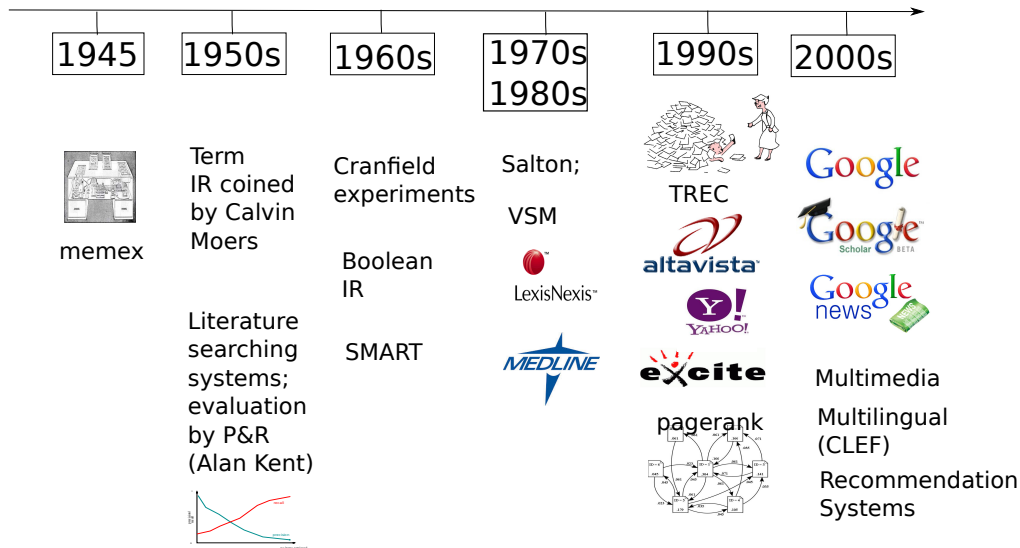


2. Structured data used in DB style searching,

   SELECT * FROM catalogue WHERE category = "florist" AND zip = "CB1"
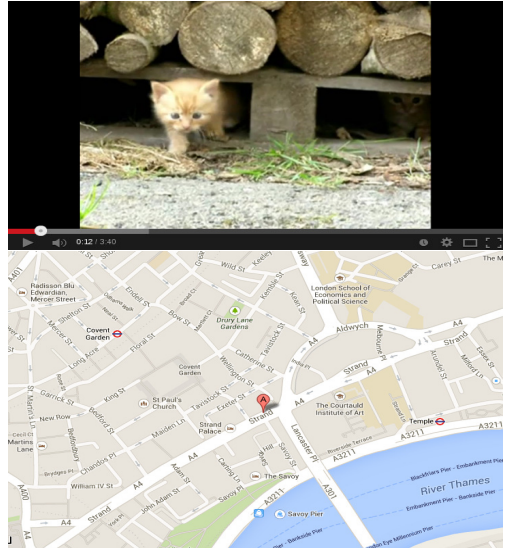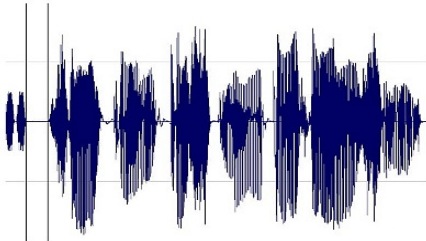
3. This does not mean that there is no structure in the data

   - Document structure (headings, paragraphs, lists. . . )

   - Explicit markup formatting (e.g. in HTML, XML. . . )

   - Linguistic structure (latent, hidden)

| 1945 | 1950s | 1960s | 1970s 1980s | 1990s | 2000s |

memex

Term IR coined by Calvin Moers

Literature searching systems; evaluation by P&R (Alan Kent)

Cranfield experiments

Boolean IR

SMART

Salton;

VSM

LexisNexis

MEDLINE

TREC

altavista

YAHOO!

excite

pagerank

Google

Google Scholar BETA

Google news

Multimedia

Multilingual (CLEF)

Recommendation Systems

# Course overview

## Course overview

1. Introduction
2. Text processing
3. Ranking with indexes
4. Retrieval models
5. Evaluation of IR systems
6. Machine Learning in IR (classification, clustering, and learning to rank)
7. Web information retrieval and search engines
8. Neural information retrieval
9. Applications
   - Recommender systems
   - Personalized IR
   - Sentiment Analysis
   - Corss-lingual IR
   - QA systems

# References

1. Chapter 1 of Information Retrieval[1]

2. Chapter 1 of Search Engines - Information Retrieval in Practice.[2]

---

[1] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze (2008). *Introduction to Information Retrieval*. New York, NY, USA: Cambridge University Press.

[2] W. Bruce Croft, Donald Metzler, and Trevor Strohman (2009). *Search Engines - Information Retrieval in Practice*. Pearson Education.

**Questions?**