# CMPE 493
# INTRODUCTION TO
# INFORMATION RETRIEVAL

## Introduction

Arzucan Özgür

Department of Computer Engineering, Boğaziçi University
September 29, 2015

---

## Course Staff

▸ Instructor: Arzucan Özgür
  ▸ Office: ETA 18
  ▸ Phone: 0212-359-7226
  ▸ E-mail: arzucan.ozgur@boun.edu.tr

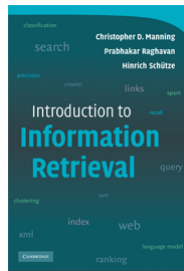  (Please include CMPE493 in your subject when sending e-mail.)
  ▸ Office hours: Monday 14:00-15:00, Tuesday 13:00-15:00, or by appointment.

▸ TAs:

▸ Şaziye Betül Özateş (sbetulbilgin@gmail.com)

▸ Alper Çetiner (alper.cetiner@boun.edu.tr)

# Text book

▸ Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval, Cambridge University Press. 2008.
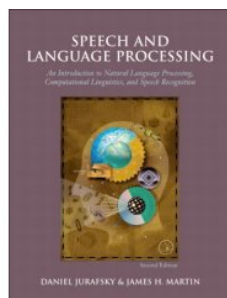


Available online (free) at the website of the book:
http://nlp.stanford.edu/IR-book/information-retrieval-book.html
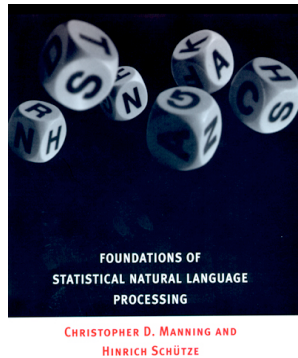
# Reference book (Optional)

▸ Daniel Jurafsky and James H. Martin, SPEECH and LANGUAGE PROCESSING: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition, Second Edition, 2008.



Available at the Bookstore.

## Reference book (Optional)

▸ Christopher D. Manning and Hinrich Schütze, Foundations of Statistical Natural Language Processing, MIT Press, 1999. http://nlp.stanford.edu/fsnlp/

**FOUNDATIONS OF STATISTICAL NATURAL LANGUAGE PROCESSING**

CHRISTOPHER D. MANNING AND HINRICH SCHÜTZE

## Course Web Site:

• We will use the Moodle Course Management System for lecture notes, announcements, homework/project submissions, and grading.

   • https://moodle.boun.edu.tr

You will automatically be subscribed to the system. You can login using your "boun" e-mail account's username and password.

## Grading

- Midterm Exam: 15%
- Final Exam: 15%
- 3-4 Assignments: 30%
- Term Project: 35%
- Class Participation: 5%

## Grading – Exams

- In-class midterm and final exams
- Consisting of problems covering the lecture material
- Closed book/notes
- Dates:
  - Midterm Exam: November 4, in the lecture hour (15:00-17:00)
  - Final Exam: As scheduled by the registration office

## Grading – Homework Assignments

- Involve some programming where you will implement and test some of the techniques that we cover in class.

- You can use any programming language of your choice such as Perl, Python, Java, etc.

- We should be able to run your program.

- You should provide a readme file, explaining how to run your program.

## Term Project

One of the aspects of this course is preparing you for original research in IR.
- ▸ Identifying an interesting problem
- ▸ Gathering relevant literature and datasets
- ▸ Solving it using new algorithms
- ▸ Evaluating the results

▸ Ability to present your ideas and research
- ▸ Writing up your results in a scientific paper format
- ▸ Presenting a research talk to a scientific audience

# Term Project

▸ The project teams can consist of one or two people. (Teams consisting of two people is recommended)

▸ Each team will choose a project topic by selecting a recent scientific paper from an IR/NLP conference or journal.

▸ The project will involve replicating the work done in the paper and proposing extensions/improvements to the existing work. The proposed extensions do not need to be implemented.

# Some of the Relevant Scientific Conferences

▸ ACM SIGIR Conference on Research and Development in Information Retrieval

▸ Conference on Information and Knowledge Management (CIKM)

▸ ACM International Conference on Web Search and Web Data Mining (WSDM)

▸ Association for Computational Linguistics (ACL)

▸ North American Association for Computational Linguistics (NAACL)

▸ Empirical Methods in Natural Language Processing (EMNLP)

▸ International Conference on Computational Linguistics (COLING)

▸ You can select your papers from relevant journals as well, including Information Retrieval, Computational Linguistics, TACL, Natural Language Engineering, and Journal of the Association for Information Science and Technology (JASIST)

## Term Project - Deliverables

‣ Paper selection and 1-2 page description of the methodology planned to be used to replicate the work
  ‣ Date: November 2

‣ Short project presentation in the end of the semester
  ‣ Tentative Dates: December 15-16, December 22-23, lecture hours

‣ Submit a project report in the end of the semester.
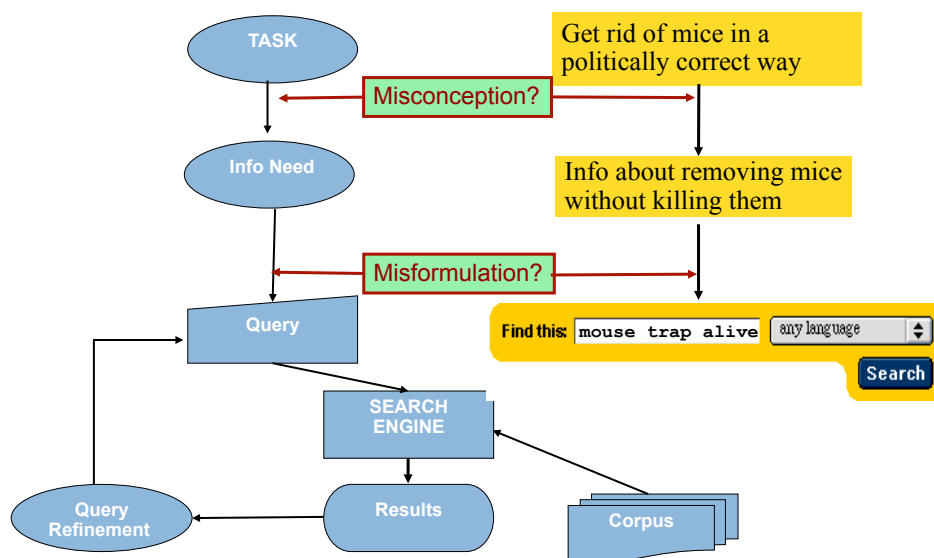  ‣ Tentative Date: Final exam date

## 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.

# Basic assumptions of Information Retrieval

▸ Collection: Fixed set of documents

▸ Goal: Retrieve documents with information that is
  <u>relevant</u> to the user's information need and helps the user
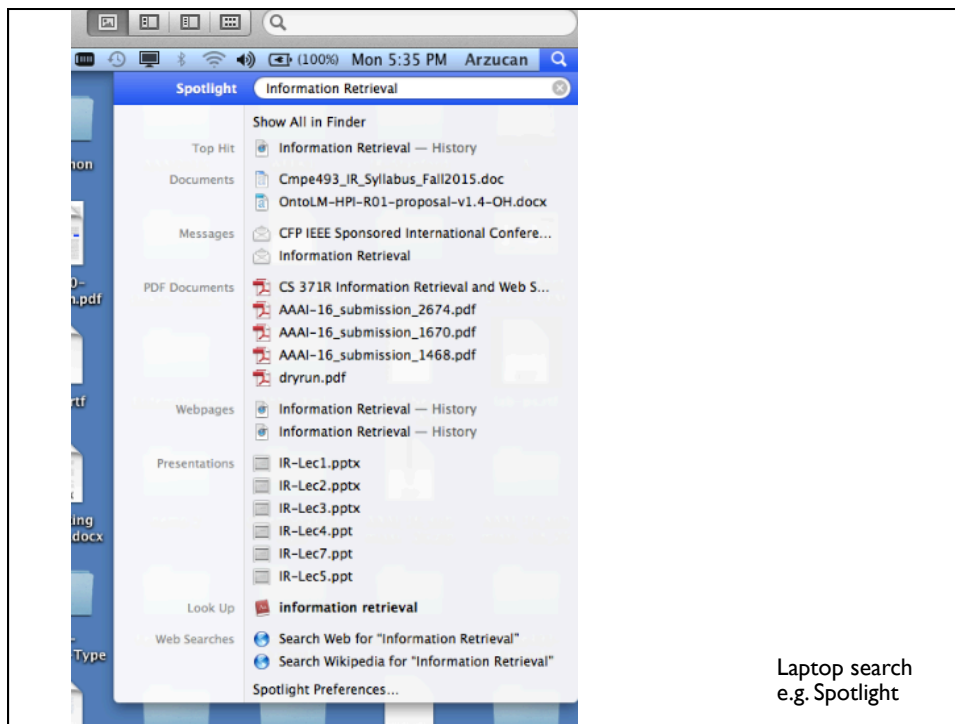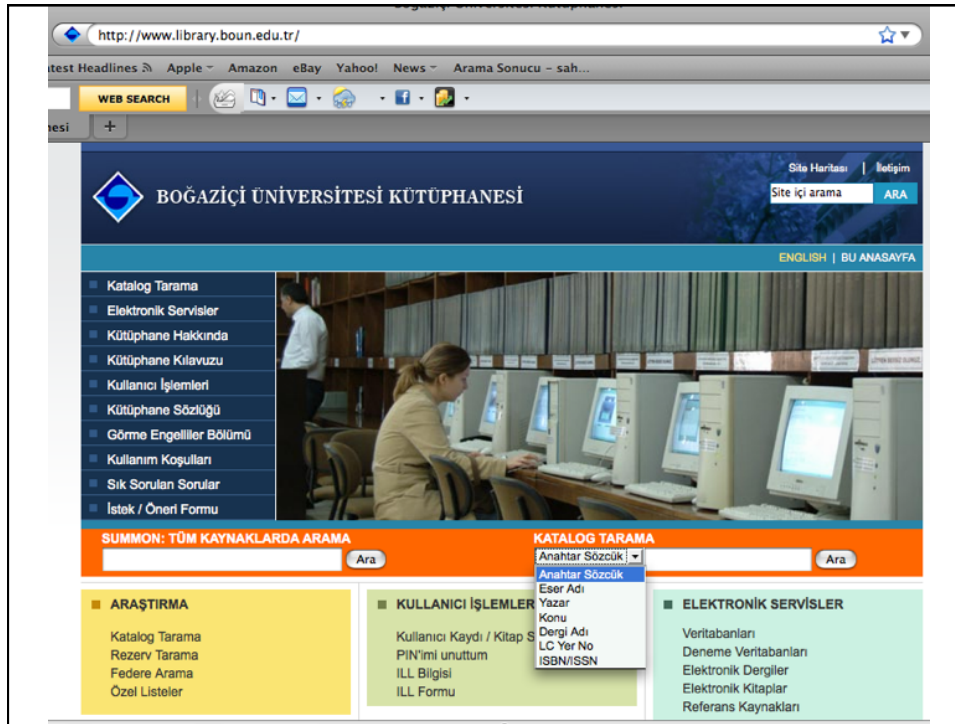  complete a task

---

# The classic search model

# How good are the retrieved docs?

▸ *Precision*: Fraction of retrieved docs that are relevant to user's information need

▸ *Recall*: Fraction of relevant docs in collection that are retrieved

▸ More precise definitions and measurements to follow in later lectures
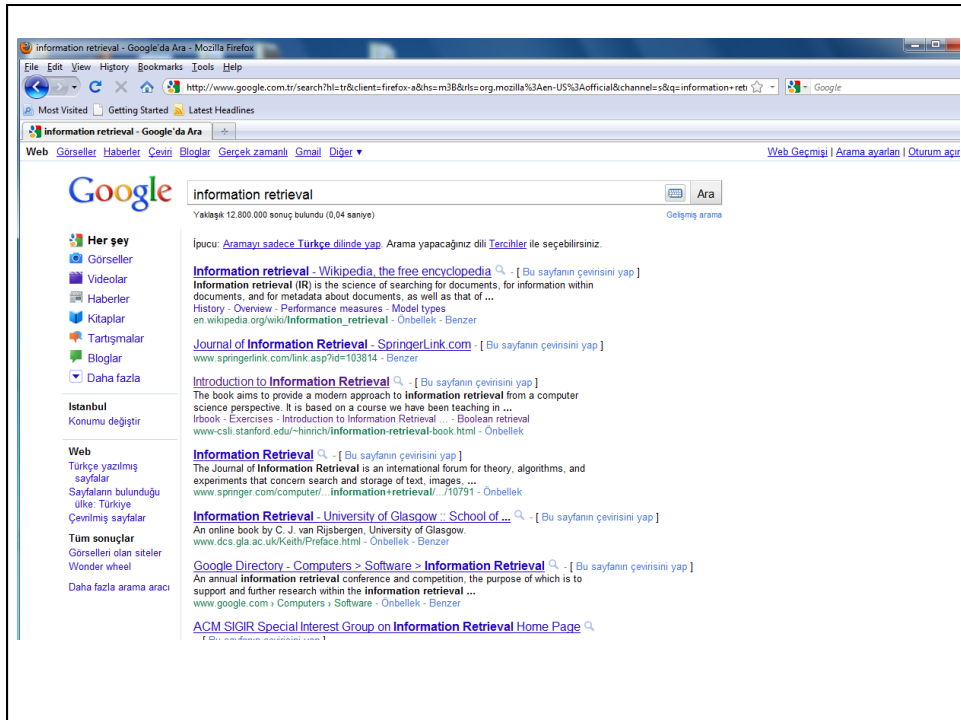
# Examples of search engines

▸ Conventional (library catalog).

Search by keyword, title, author, etc.

▸ Text-based (Google, Yahoo!, Bing, Yandex, Baidu; also email search, laptop search etc.)

Search by keywords. Limited search using queries in natural language.

▸ Multimedia (QBIC, WebSeek)

Search by visual appearance (shapes, colors,… ).

▸ Question answering systems (Ask, NSIR, Answerbus)

Search in (restricted) natural language

▸ Other:

music retrieval

Laptop search
e.g. Spotlight

# IR systems on the Web

▸ Search for Web pages: http://www.google.com

▸ Domain specific search (e.g., legal, biomedical): PubMed

▸ Search for images: http://www.picsearch.com

▸ Search for image content: http://wang14.ist.psu.edu/

▸ Search for answers to questions: http://www.ask.com

▸ Music retrieval: http://www.rotorbrain.com/foote/musicr/

**Screenshot 1 — Google News (Mozilla Firefox)**

european union - Google News - Mozilla Firefox

File  Edit  View  History  Bookmarks  ScrapBook  Tools  Help

www.google.com

Sign in

Web   Images   Video   News   Maps   more »

Google News

european union

Search News    Search the Web

Advanced news search
Preferences

Results **1 - 10** of about 27,833 for **european union**. (0.71 seconds)

Try your search on: Yahoo News, Ask, AllTheWeb, MSN, Lycos, Sky News, CNN, Feedster, Daypop, Bloglines

Sorted by relevance    Sort by date

Top Stories
World
U.S.
Business
Sci/Tech
Sports
Entertainment
Health
Most Popular

News Alerts

RSS | Atom
About Feeds

Mobile News

About
Google News

**EXTREME SOLIDARITY Far-Right Parties Form New Group in European ...**
Spiegel Online, Germany - 42 minutes ago
**European Union** expansion is a topic typically supported by those on the left of the continent's political spectrum and opposed by those on the right. ...
Far-right EU lawmakers form coalition  Olberlin
all 88 news articles »

**Wild bird trade to be banned by European Union**
EnjoyFrance.com, France - 1 hour ago
The **European Union** is going to ban the trade in wild birds starting in July, EU animal health officials have announced. Animal welfare campaigners are ...
Wild bird imports to end  Green Consumer Guide
UN-Backed Body 'Disappointed' By Bird Trade Ban  Scoop.co.nz (press release)
EU To Ban Wild Birds Imports  All Headline News
Earthtimes.org
all 8 news articles »

Severe Weather Alert

**Russia, European Union A serious problem of trust**
Monday Morning, Lebanon - 5 hours ago
Merkel, in contrast, is wary of depending heavily on Russia for oil and gas and

Find:        Next  Previous  Highlight all

Done    FoxyT    Open Notebook

---

**Screenshot 2 — Ask.com (Mozilla Firefox)**

Ask.com - What's Your Question? - Mozilla Firefox

File  Edit  View  History  Bookmarks  Tools  Help

http://www.ask.com/web?qsrc=2990&o=10181&l=dir&q=What+is+the+capital+of+Turkey%3F

Google

Most Visited    Getting Started    Latest Headlines

Ask.com - What's Your Question?

Community    Web    Images    News    Videos    More »    Advanced Search    Settings    Sign In

Ask    What is the capital of Turkey?    Search

Top Answer

The Capital of Turkey is Ankara.

Source: CIA World Factbook
See Also: BBC Profile · Encyclopedia
Search For: Flights · Geography · Government · People

Source

Was this answer helpful?

Community    122623 people answering

What is the capital of Turkey?

Ask the Community ›

New from Ask. See how it works »

**Work and Travel**    Ads
Tecrübe Güven ve Kalite'de 11. Yıl -Kariyer Programlarında Bir Marka
www.armadagrandee.com

**Airport Hotel ISTANBUL**
Size zaman kazandırmak için tasarlandı...
www.isgairporthotel.com

**Made-in-Turkey**
Türk Üreticileri ve Sanayicileri İhracat için sanal mağazanız burada
www.made-in-turkey.com

**İstanbul İş İlanları**
İstanbul'un En İyi Firmalarında İş İmkanı Monster'da Hemen Üye Ol
monster.com.tr/İstanbul-is-ilanlari

**Related Searches**
Map of Turkey
Map of Europe
Map of Middle East
Turkey Country Information
History of Istanbul
Ankara
World Atlas
Constantinople
Map of Spain
Cyprus
Map of Italy
Map of Africa

**What is the capital of turkey?**
The capital of Turkey is Ankara. Turkey is a country located in the Middle East. They are the US's closest ally in that region, other than Israel.
http://answers.ask.com/Society/Government_and_Law/what_...

**Related Questions**

Türkiye'nin coğrafi bölgeleri nelerdir?

Sor

Yabancı Kaynaklar

Sucuk döneri: Afyonkarahisar mutfağına özgü lezzetli sucuktan yapılan döner türüdür.

Türkiye'de İslam en yaygın dindir.

Türkiye'nin coğrafi bölgeleri, 6 Haziran - 21 Haziran 1941 tarihleri arasında Ankara'da toplanan Birinci Coğrafya Kongresi tarafından belirlenmiştir. Bu çalışmanın sonucunda Türkiye'nin üç tarafının denizlerle çevrilmiş olması, dağların Anadolu'nun iç kesimlerini kıyılardan ayırması, iklim, ulaşım ve bitki örtüsü gibi kriterler dikkate alınarak Türkiye'nin coğrafi bölgeleri belirlenmiştir.
Coğrafi bölgeleri oluşturan etkenler.
Coğrafi bölgeler ve coğrafi bölgelerin sınırları belirlenirken şu etkenler dikkate alınmıştır;
Bölgeler ve bölümler.
Doğal, beşerî ve ekonomik özellikler yönünden sınırları içinde benzerlik gösteren geniş alanlara bölge denir.
Sınırları içinde benzerlikleri olan ancak bölgenin diğer yerlerinden farklı olan küçük alanlara ise bölüm denir.
Birinci Coğrafya Kongresinde Türkiye coğrafi 7 bölgeye ve 21 bölüme ayrılmıştır.
Türkiye'nin yedi coğrafi bölgesinden dördüne komşu olduğu denizin adı verilmiştir (Akdeniz Bölgesi, Karadeniz Bölgesi, Ege Bölgesi, Marmara Bölgesi).
Diğer üç bölge de Anadolu bütünü içindeki konumlarına göre adlandırılmışlardır (İç Anadolu Bölgesi, Doğu Anadolu Bölgesi, Güneydoğu Anadolu Bölgesi).
Türkiye'deki coğrafi bölgeler arasında nüfus miktarı ve yoğunluğu yönünden önemli farklar bulunmaktadır.
Nüfusun en yoğun olduğu bölge Marmara Bölgesi en seyrek olduğu bölge de Doğu

C. Derici, T. Güngör, et al.,
Question Analysis for a Closed Domain
Question Answering System, CICLing 2015

## What does it take to build a search engine?

- Decide what to index
- Collect it
- Index it (efficiently)
- Keep the index up to date
- Provide user-friendly query facilities

## What else?

- Understand the structure of the web for efficient crawling
- Understand user information needs
- Preprocess unstructured textual data
- Cluster data
- Classify data
- Evaluate performance

## Goals of the course

▸ Understand how search engines work

▸ Understand the limits of existing search technology

▸ Learn about the state of the art in IR research

▸ Learn to analyze textual data sets

▸ Learn to evaluate information retrieval systems

▸ Learn about standardized document collections

▸ Learn about text similarity measures

▸ Learn about semantic dimensionality reduction

▸ Learn about web crawling

▸ Learn to use existing software

▸ Understand the dynamics of the Web by building appropriate mathematical models

▸ Build working systems that assist users in finding useful information from large collections


## Topics (tentative list)

▸ Boolean model; text pre-processing; inverted indexes

▸ Approximate string matching and tolerant retrieval

▸ Index construction and compression

▸ Vector space model; text-similarity metrics; term weighting; ranked retrieval

▸ Evaluating information retrieval systems

▸ Relevance feedback; query expansion

▸ Language models for information retrieval

▸ Text classification and clustering

▸ Latent semantic indexing

▸ Web search and crawling

▸ Link analysis (e.g. hubs and authorities, Google PageRank)

## References

▸ Content adapted from Prof. Dragomir Radev and the IR book's web site.