# Introduction to Information Retrieval

Hongning Wang

CS@UVa

# What is information retrieval?

# What is information retrieval?

- Apple's vision 35 years ago

[Knowledge Navigator](Knowledge Navigator)

# Why information retrieval

- Information overload
  - *"It refers to the <u>difficulty</u> a person can have understanding an issue and making decisions that can be caused by the presence of <u>too much</u> information."* - wiki

# Why information retrieval

- Information overload



Hobbes' Internet Timeline Copyright ©2014 Robert H Zakon
http://www.zakon.org/robert/internet/timeline/

| DATE | SITES | | DATE | SITES |
|------|-------|---|------|-------|
| 12/90 | 1 | | 06/95 | 23,500 |
| 12/91 | 10 | | 01/96 | 100,000 |
| 12/92 | 50 | | 06/96 | 252,000 |
| 06/93 | 130 | | 01/97 | 646,162 |
| 09/93 | 204 | | 06/97 | 1,117,259 |
| 10/93 | 228 | | 01/98 | 1,834,710 |
| 12/93 | 623 | | 06/98 | 2,410,067 |
| 06/94 | 2,738 | | 01/99 | 4,062,280 |
| 12/94 | 10,022 | | 07/99 | 6,598,697 |

Figure 2: Growth of WWW

Figure 1: Growth of Internet

# Why information retrieval

- Handling <u>unstructured</u> data
  - Structured data: database system is a good choice
  - Unst...
    - Te... lio, video...
    - "s... as ur...
    - U



Table 1: People in CS Department

| ID | Name | Job |
|----|------|-----|
| 1 | Jack | Professor |
| 3 | David | Stuff |
| 5 | Tony | IT support |

select Name from People where Job = 'Stuff';

Total Enterprise Data Growth 2005-2015, IDC 2012

# Why information retrieval

- An essential tool to deal with information overload

You are here!

# History of information retrieval

- Idea popularized in the pioneer article "***As We May Think***" by Vannevar Bush, 1945
  - *"Wholly new forms of <u>encyclopedias</u> will appear, ready-made with a mesh of <u>associative trails</u> running through them, ready to be dropped into the memex and there amplified."* **-> WWW**
  - *"A memex is a device in which an individual <u>stores all</u> his books, records, and communications, and which is mechanized so that it may be consulted with <u>exceeding speed</u> and <u>flexibility</u>."* **-> Search engine**

# Major research milestones

- Early days (late 1950s to 1960s): foundation of the field
  - Luhn's work on automatic indexing
  - Cleverdon's Cranfield evaluation methodology and index experiments
  - Salton's early work on SMART system and experiments
- 1970s-1980s: a large number of retrieval models
  - Vector space model
  - Probabilistic models
- 1990s: further development of retrieval models and new tasks
  - Language models
  - TREC evaluation
  - Web search
- 2000s-present: more applications, especially Web search and interactions with other fields
  - Learning to rank
  - Scalability (e.g., MapReduce)
  - Real-time search

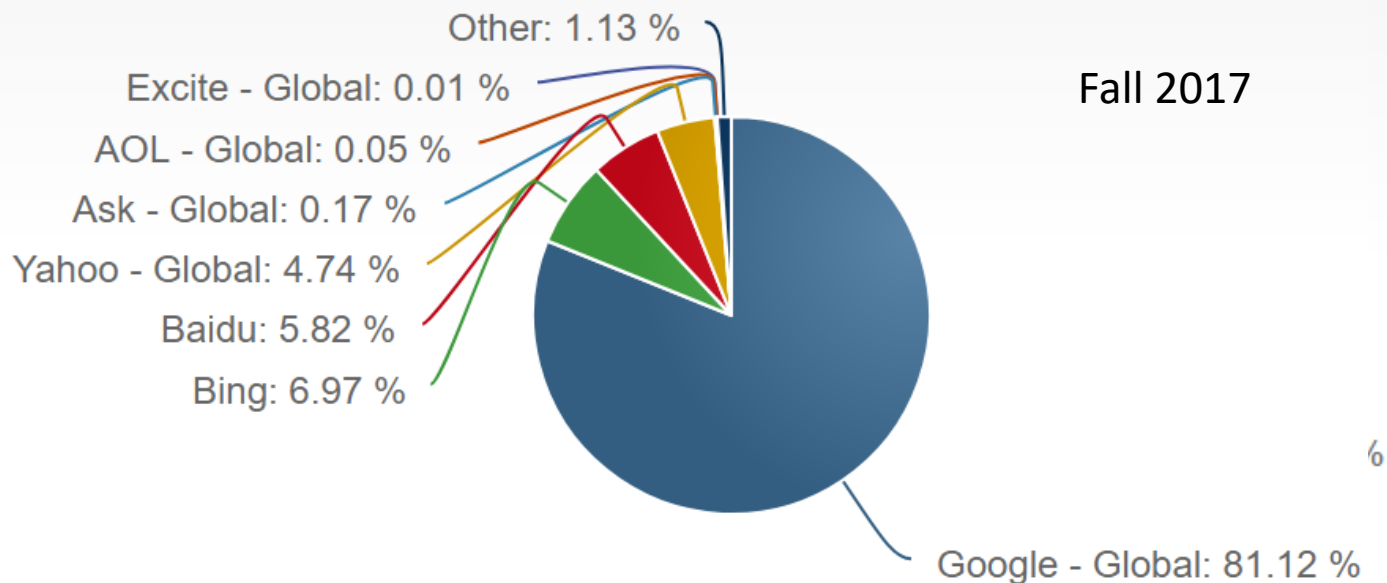# History of information retrieval

- Catalyst
  - Academia: Text Retrieval Conference (TREC) in 1992
    - *"Its purpose was to support research within the information retrieval community by providing the infrastructure necessary for large-scale <u>evaluation</u> of text retrieval methodologies."*
    - *"… about <u>one-third</u> of the improvement in web search engines from 1999 to 2009 is attributable to TREC. Those enhancements likely saved up to <u>3 billion hours</u> of time using web search engines."*
    - Till today, it is still a major test-bed for academic research in IR

# History of information retrieval

- Catalyst
  - Industry: web search engines
    - WWW unleashed explosion of published information and drove the innovation of IR techniques
    - First web search engine: "*Oscar Nierstrasz at the University of Geneva wrote a series of Perl scripts that <u>periodically mirrored</u> these pages and rewrote them into a <u>standard format</u>*." Sept 2, 1993
    - Lycos (started at CMU) was launched and became a major commercial endeavor in 1994
    - Booming of search engine industry: *Magellan, Excite, Infoseek, Inktomi, Northern Light, AltaVista, Yahoo!, Google,* and *Bing*
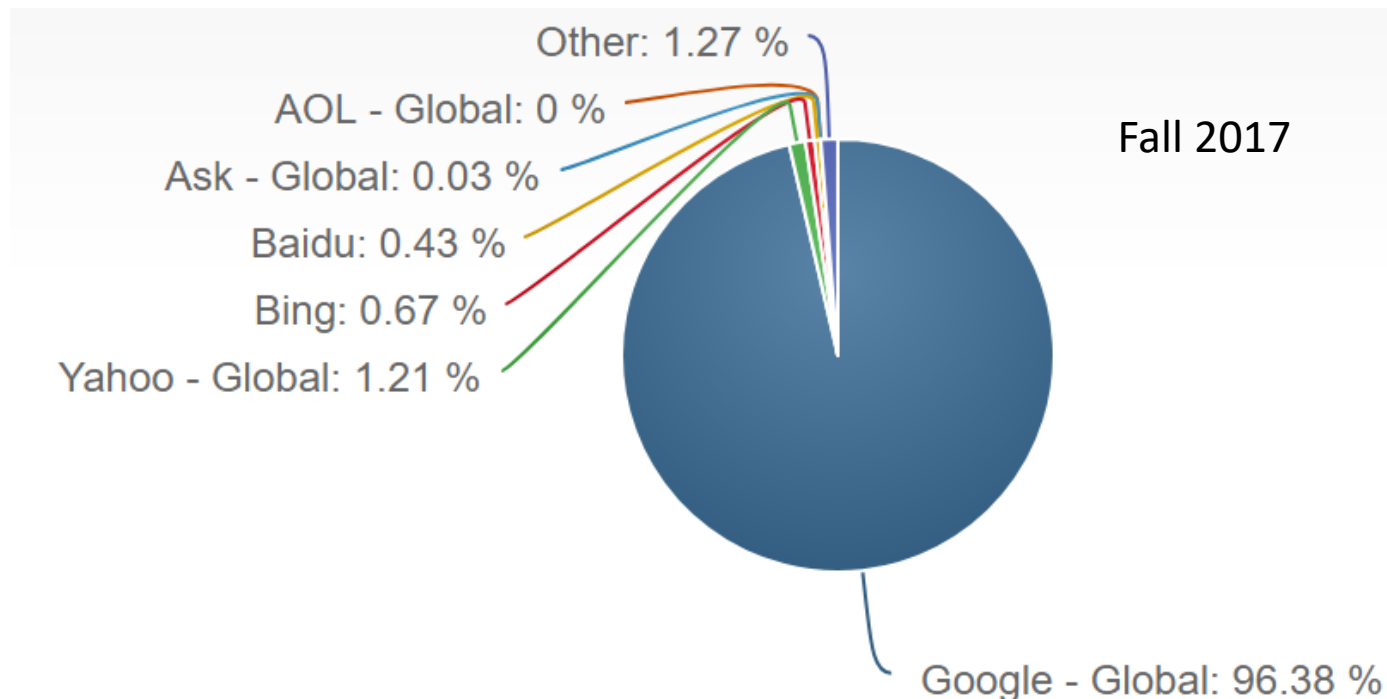
# Major players in this game

- Global search engine market - desktop
  - By http://marketshare.hitslink.com/search-engine-market-share.aspx



Other: 1.13 %
Excite - Global: 0.01 %
AOL - Global: 0.05 %
Ask - Global: 0.17 %
Yahoo - Global: 4.74 %
Baidu: 5.82 %
Bing: 6.97 %
Fall 2017
Google - Global: 81.12 %

# Major players in this game

- Global search engine market - mobile
  - By http://marketshare.hitslink.com/search-engine-market-share.aspx

Other: 1.27 %

AOL - Global: 0 %

Ask - Global: 0.03 %

Baidu: 0.43 %

Bing: 0.67 %

Yahoo - Global: 1.21 %

Fall 2017

Google - Global: 96.38 %

# How to perform information retrieval

- Information retrieval when we did not have a computer

# Welcome back

- We will start our discussion at 2pm

- sli.do event code: <u>32233</u>

- TA's office hour has been finalized at <mark>Mon/Wed, 1pm-2pm</mark>, via zoom, by appointment

# Recap: why information retrieval

- ## Handling <u>unstructured</u> data
  - – Structured data: database system is a good choice
  - – Unstructured data is more dominant
    - Text in Web documents or emails, image, audio, video…
    - "*<u>85 percent</u> of all business information exists as unstructured data*" - Merrill Lynch
    - Unknown <u>semantic</u> meaning

# How to perform information retrieval



Crawler and indexer

Query parser

Document analyzer

Ranking model

# Crack into Google!

# Crack into Google!

Object Detection

Joseph

Result Ranker

(Jay)

Personal Preferences (Henry)

Page Crawler (Sammy)

remove stopwords, stemwords

Search Decoder amar

repository

# How to perform information retrieval



**PARSING & INDEXING**

Doc Rep

Query Rep

query

**Repository**

**User**

Ranking

**SEARCH**

**APPLICATIONS**

results

**LEARNING**

Evaluation

judgments

**FEEDBACK**

**We will cover:**

1) Search engine architecture;  2) Retrieval models;

3) Retrieval evaluation;  4) Relevance feedback;

5) Link analysis;  6) Search applications.

# Core concepts in IR

- Query representation
  - Lexical gap: say v.s. said
  - Semantic gap: ranking model v.s. retrieval method
- Document representation
  - Special data structure for efficient access
  - Lexical gap and semantic gap
- Retrieval model
  - Algorithms that find the ***most relevant*** documents for the given information need

# A glance of modern search engine

- In old times

*Yet Another **Hierarchical** Officious/Obstreperous/ Odiferous/Organized **Oracle***



*Yahoo race of fictional beings from **Gulliver's Travels***

# A glance of modern search engine



Demand of understanding

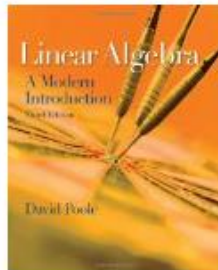Demand of efficiency

Demand of convenience

Demand of diversity

# IR is not just about web search

- Web search is just one important area of information retrieval, but not all
- Information retrieval also includes
  - Recommendation



Recommended Based on Your Browsing History

Linear Algebra and Its Applications...
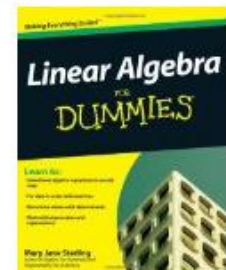> David C. Lay
Hardcover
★★★☆☆ (84)
$183.33 $141.16

Linear Algebra: A Modern Introduction
> David Poole
Hardcover
★★★★☆ (41)
$316.95 $289.88

Linear Algebra
> G. E. Shilov
Paperback
★★★★☆ (34)
$18.95 $12.65

Introduction to Linear Algebra...
> Gilbert Strang
Hardcover
★★★☆☆ (57)
$87.50 $83.13

Linear Algebra For Dummies
> Mary Jane Sterling
Paperback
★★★★☆ (29)
$19.99 $16.23

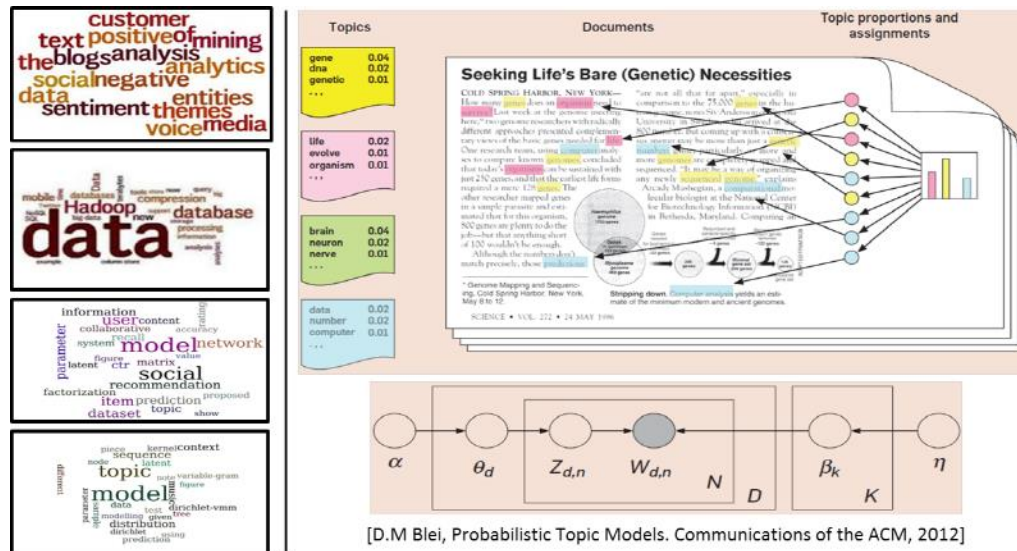# IR is not just about web search

# IR is not just about web search

- Web search is just one important area of information retrieval, but not all

- Information retrieval also includes
  - Question answering

# IR is not just about web search

- Web search is just one important area of information retrieval, but not all

- Information retrieval also includes
  - Text mining



[D.M Blei, Probabilistic Topic Models. Communications of the ACM, 2012]

# IR is not just about web search

- Web search is just one important area of information retrieval, but not all
- Information retrieval also includes
  - Online advertising

# IR is not just about web search

- Web search is just one important area of information retrieval, but not all

- Information retrieval also includes
  - Enterprise search: web search + desktop search

CS 4780: Information Retrieval

# Related Areas



Mathematics

Applications

Statistics
Optimization

Machine Learning
Pattern Recognition

Natural
Language
Processing

Web Applications,
Bioinformatics…

**Information
Retrieval**

Library & Info
Science

Databases

Data Mining

Software engineering
Computer systems

Algorithms

Systems

# IR v.s. DBs

- Information Retrieval:
  - Unstructured data
  - Semantics of objects are subjective
  - Simple keyword queries
  - Relevance-drive retrieval
  - Effectiveness is primary issue, though efficiency is also important

- Database Systems:
  - Structured data
  - Semantics of each object are well defined
  - Structured query languages (e.g., SQL)
  - Exact retrieval
  - Emphasis on efficiency

# IR and DBs are getting closer

- IR => DBs
  - Approximate search is available in DBs
  - Eg. in mySQL

  **mysql> SELECT * FROM articles
    -> WHERE MATCH (title,body)
    AGAINST ('database');**

- DBs => IR
  - Use information extraction to convert unstructured data to structured data, e.g., knowledge base
  - Semi-structured representation: XML data; queries with structured information

# IR v.s. NLP

- Information retrieval
  - Computational approaches
  - Statistical (shallow) understanding of language
  - Handle large scale problems

- Natural language processing
  - Cognitive, symbolic and computational approaches
  - Semantic (deep) understanding of language
  - (often times) small scale problems

# IR and NLP are getting closer

- IR => NLP
  - Larger data collections
  - Scalable/robust NLP techniques, e.g., translation models

- NLP => IR
  - Deep analysis of text documents and queries
  - Information extraction for structured IR tasks
  - Natural language based QA systems

# Text books



- ***Introduction to Information Retrieval***. Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schuetze, Cambridge University Press, 2007.



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

# What to read?

**Mathematics**

**Web Applications, Bioinformatics…**

**Machine Learning Pattern Recognition**
**ICML, NIPS, UAI**

**Library & Info Science**

**Statistics Optimization**
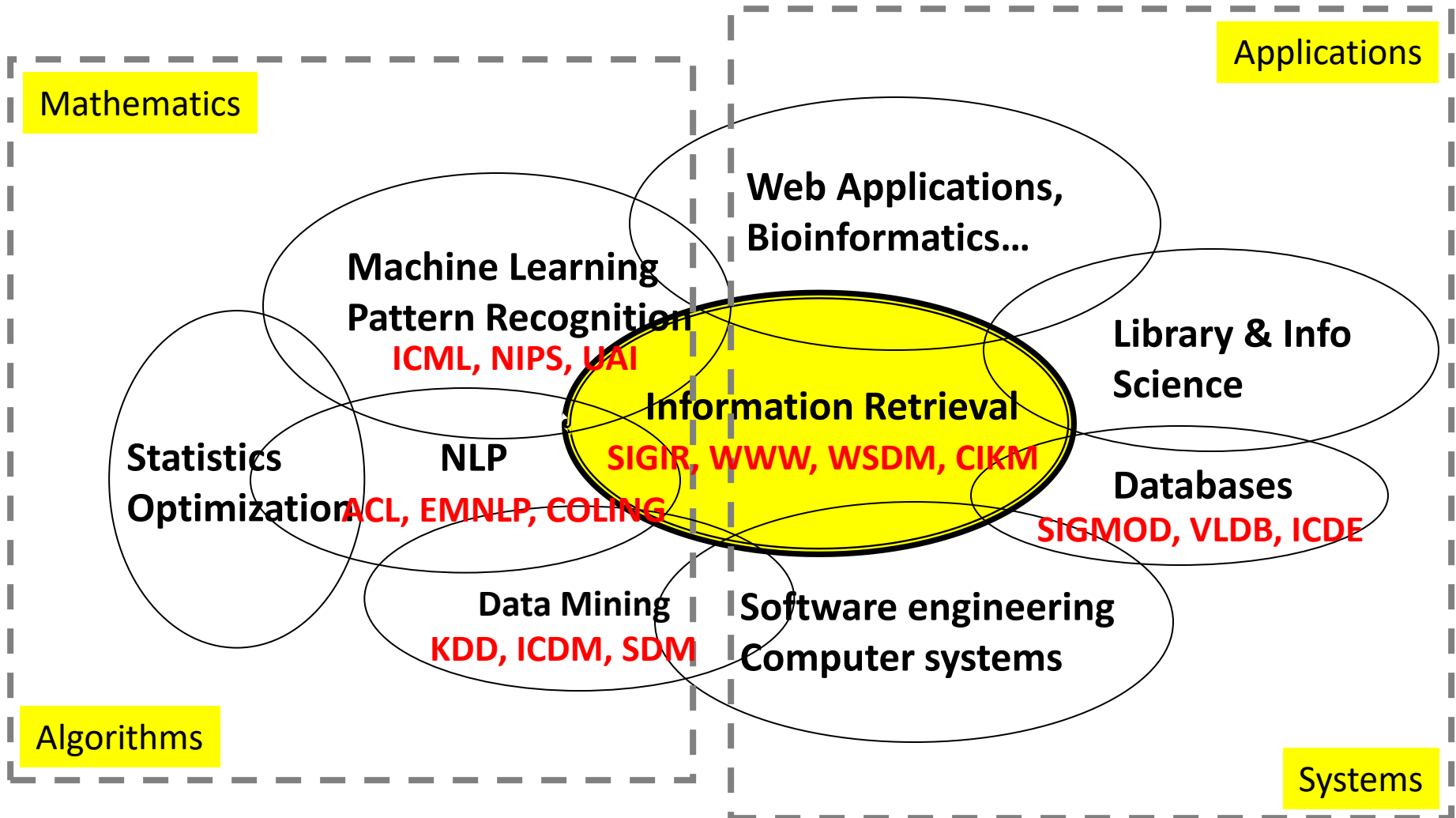**Information Retrieval**
**SIGIR, WWW, WSDM, CIKM**

**NLP**
**ACL, EMNLP, COLING**

**Databases**
**SIGMOD, VLDB, ICDE**

**Data Mining**
**KDD, ICDM, SDM**

**Software engineering Computer systems**

**Algorithms**

**Systems**

- Find more on course website for resource

# IR in future

- Mobile search
  - Desktop search + location? Not exactly!!
- Interactive retrieval
  - Machine collaborates with human for information access
- Personal assistant
  - Proactive information retrieval
  - [Knowledge navigator](#)
- And many more
  - You name it!

# What you should know

- IR originates from library science for handling <u>unstructured</u> data

- IR has many important application areas, e.g., web search, recommendation, and question answering

- IR is a highly interdisciplinary area with DBs, NLP, ML, HCI

# Today's reading

- *Bush, Vannevar. "As we may think." The atlantic monthly 176, no.1 (1945): 101-108.*

- Introduction to Information Retrieval
  - Chapter 1: Boolean Retrieval