# Introduction to Information Retrieval

Lecture 1

CS 410/510

Information Retrieval on the Internet

# Lecture Outline

- IR systems
  - Overview
  - IR systems vs. DBMS
  - Types, facets of interest
- User tasks
- Document representations
- Queries
- Retrieval
- Evaluation

## Information retrieval

"Information retrieval (IR) deals with the representation, storage, organization of, and access to information items."

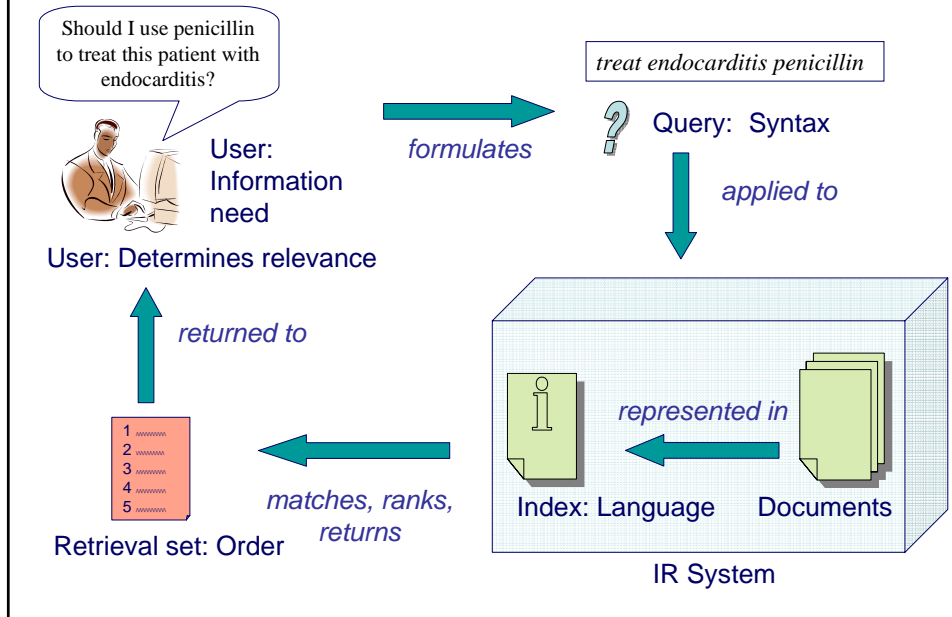- Baeza-Yates and Berthier Ribeiro-Neto in <u>Modern Information Retrieval</u>, p. 1

"*Information retrieval* is often regarded as being synonymous with *document retrieval* and, nowadays, with *text retrieval*, implying that the task of an IR system is to retrieve documents or texts with information content that is *relevant* to a user's information *need* ... the approaches that have been developed for this purpose are also applicable to a whole family of related information processing tasks that lie between, on the one hand, *data retrieval* and, on the other, *fact* or *knowledge retrieval*."

- Sparck Jones and Willett in <u>Readings in Information Retrieval</u>, p. 1

# Information Retrieval

- History
  - Early work in IR in the 50's and 60's
  - Roots in library science much older, e.g.
    - Dewey Decimal system 1870s
    - Library of Congress Classification 1890s
- Important Components:
  - Indexing
  - Searching
  - User-system interaction
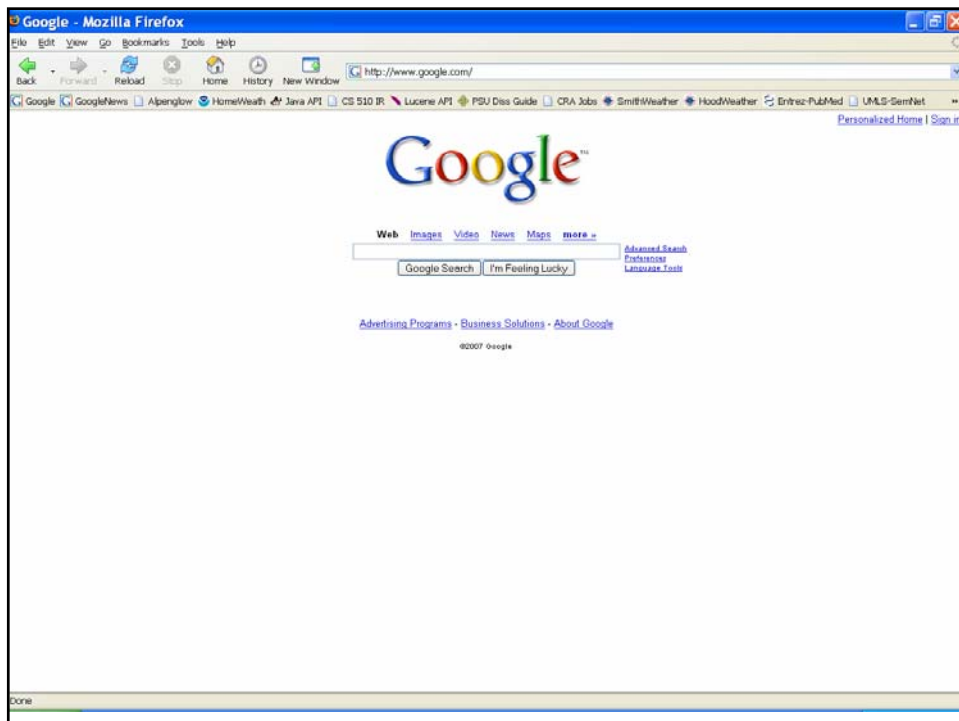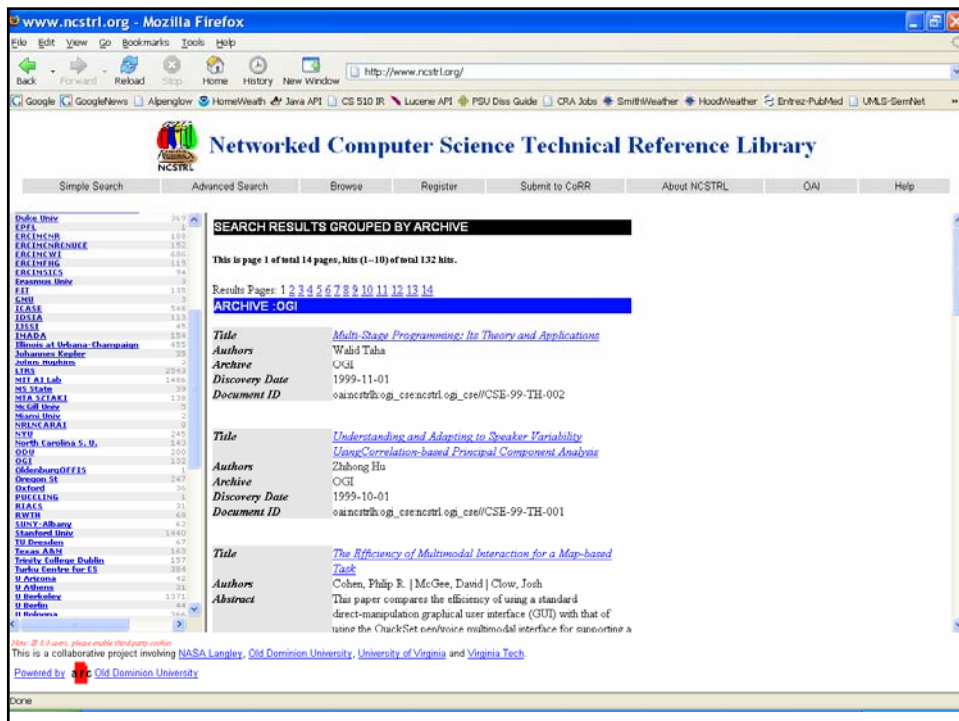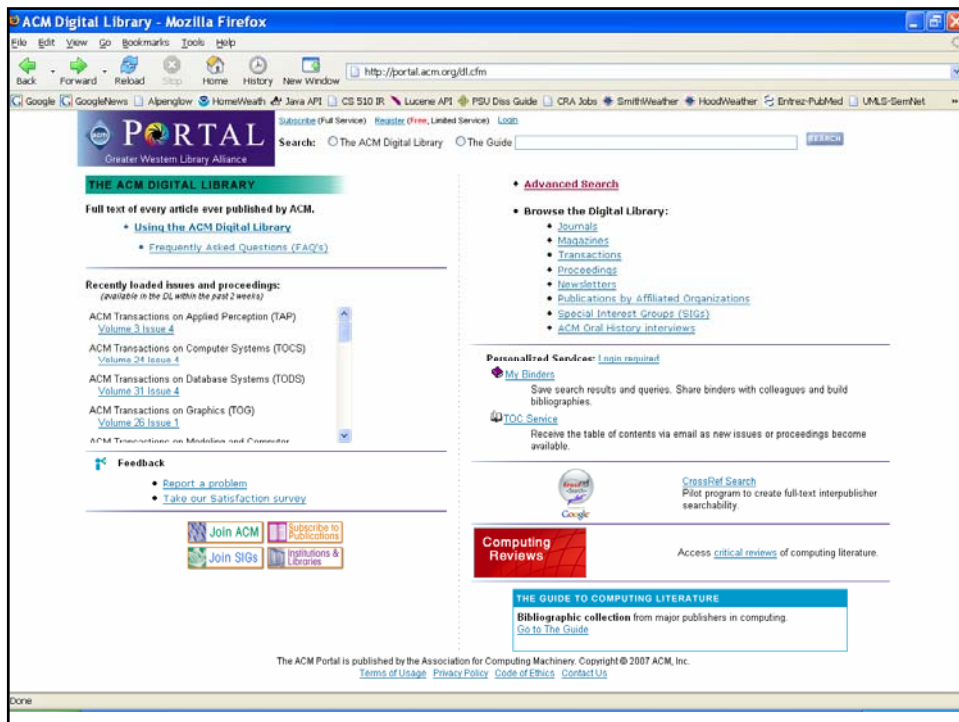
## Typical information retrieval (IR) system

Should I use penicillin to treat this patient with endocarditis?

User: Information need

User: Determines relevance

*formulates*

*treat endocarditis penicillin*

Query: Syntax

*applied to*

*returned to*

1 ......
2 ......
3 ......
4 ......
5 ......

Retrieval set: Order

*matches, ranks, returns*

*represented in*

Index: Language    Documents

IR System

---

# IR systems vs. DBMS

|  | DBMS | IR System |
|---|---|---|
| Target |  |  |
| Queries / Query language |  |  |
| Matching |  |  |
| Results |  |  |

# General types of IR systems

- Web
- Full text documents
- Bibliographic
- Distributed variations
  - Metasearch
  - Virtual document collections

6

# Facets of interest in IR systems

- Site
  - single controlled website (e.g. intranet)
  - distributed across multiple sites

# Facets of interest in IR systems

- Documents
  - Hyperlinked?
  - Format?
    - HTML
    - PDF
    - Word processed
    - Scanned OCR?
  - Type?
    - Text?
    - Multimedia?
    - Semistructured (XML)?
  - Static or dynamic?

# Facets of interest in IR systems

- Collection
  - Closed or open?
  - Curated?
  - Static or constantly changing?
  - Content?
    - specific domain?
    - general/any domain?
  - Purpose?
  - Audience?
  - Access? (restricted or not)

# Facets of interest in IR systems

- Search engine
  - Supported tasks: search, browse
  - Basic model: Boolean vs. ranked results
  - Indexing language: controlled keywords vs. natural language
  - Indexing target: bibliographic data vs. full text
  - Search syntax; available operations

# User Tasks

- Mode
  - Search (retrieve)
  - Browse
- Purpose
  - Overview
  - Question answering/fact finding
  - Comprehensive research
  - Finding known item (document, page, or site)
  - Transaction (e.g. buy a book, download a file)

# User Tasks

- Examples of queries to:
  - Find an overview?
  - Find a fact/answer a question?
  - Find comprehensive information?
  - Find a known item (document, page, or site)?
  - Find a site to execute a transaction (e.g. buy a book, download a file)?

# Document Representation: Logical

- Goal:
  - Represent the content
  - Represent other aspects (sometimes)
- Methods:
  - Assign descriptors (usually selected from a predetermined list)
  - Extract features (usually words or phrases if text document)
  - Other (e.g., try to anticipate what questions the document could answer)

# Document Representation: Logical

- History
  - Early systems mostly bibliographic
    - Systems contained brief surrogates for documents
      - Title, author, abstract, location of full document
      - Remember card catalogues??
    - Provided reference to the full document
  - Limited by available storage and processing power
  - Affected the way document content was represented

# Document Representation: Logical

- Methods:
  - Manual
    - usually assign terms from a controlled vocabulary
  - Automatic
    - usually extract terms from the document
- Considerations:
  - Size of the representation
  - Improve likelihood of appropriate matching to queries

# Document Representation: Logical

**Abstract:** The objectives of this study were to determine if (1) children with migraine experience greater sleep disturbances than their siblings, (2) those with more severe migraine have greater levels of sleep disruption, and (3) these sleep disturbances lead to greater behavioral problems and more missed school. Children aged 6 to 18 years with a diagnosis of migraine for > 6 months, who had at least one sibling without migraine in the same age range, were identified through our neurology clinic database or at the time of the clinic visit. Parents completed the (1) demographic, general health, and migraine information questionnaire; (2) Child Sleep Habits Questionnaire; and (3) Behavior Assessment System for Children: Second Edition (BASC-2) Parent Rating Scales for each child. Cases with migraine had higher total sleep (P < .02), sleep delay (P < .03), and daytime sleepiness scores (P < .001) than controls. Cases with more severe migraines had higher total sleep (P < .01) and sleep duration scores (P < .03) than those with milder headaches. In cases, higher total sleep ...

**Assigned Indexing Terms:**
Activities of Daily Living
Adolescent
Case-Control Studies
Child
Child Behavior
Disorders/*complications/psychology
Circadian Rhythm/physiology
Female
Humans
Male
Migraine
Disorders/*complications/psychology
Severity of Illness Index
Sleep
Disorders/*complications/psychology
Sleep Stages/physiology
Wakefulness/*physiology

# Document Representation: Logical

- Techniques to improve representation
  - Remove stopwords
  - Stemming
  - Use document structure
- Possible text processing steps
  - Determine *what* to index (frames? page title? metadata?)
  - Strip formatting (e.g. html tags, w.p. instr.)
  - Recognize structure
  - Normalize representation (Normalize what?)
  - Recognize words (or phrases)

# Document Representation: Physical

- Associate document identifier with various descriptors
- If descriptors are extracted terms, may also include
  - Frequency of terms
  - Position of terms
- Considerations
  - Efficient storage
  - Representations that can be searched to allow a fast response to user requests

# Queries

- User has an information need
- Information need communicated to IR system as a request
- Request encoded in a query expression
- Query expression must be interpretable by the system (query language)
  - Encoded in the interface
  - Query expressions that use symbols to express operators

# Retrieval

- Match information request representation (query) to document representations (stored in an index)
- Use an algorithm to compute the matching
  - Exact match (Boolean expressions)
    - How should results be ordered?
  - Ranked list based on similarity computation (and possibly other factors)
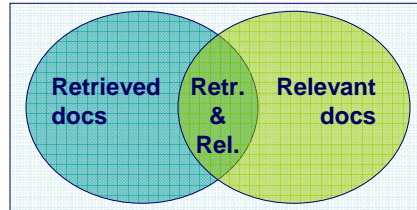
# IR as Classification

- Indexing techniques classify documents
  - Keywords, either assigned or extracted, define categories
- Information requests also identify one or more categories of interest
- Compare to bricks-and-mortar libraries?

# Evaluation

- How can we evaluate performance of an IR system?
  - System perspective
  - User perspective
- Relevance
  - (How well) does a document satisfy the user's need (or "anomalous state of knowledge")?

# Evaluation: Metrics

- There are many
- For now, two basics:



Retrieved docs | Retr. & Rel. | Relevant docs

$$\text{Recall} = \frac{\text{\# documents retrieved and relevant}}{\text{\# documents relevant}}$$

$$\text{Precision} = \frac{\text{\# documents retrieved and relevant}}{\text{\# documents retrieved}}$$

---

# Questions?

- Next: Overview of text processing