# Information Retrieval

# Information Retrieval

# Domains, Applications, and Tasks

- Web search
- Vertical search
- Enterprise search
- Media search
- Question answering
- Recommender systems
- Advertising
- Personal item search
- Passage retrieval

- Filtering
- Summarization
- Clustering
- Topic detection
- Cross-language
- Federated search
- Metasearch
- Social search
- Novel-item retrieval

# What is IR?

- Gerard Salton, 1968:
  - *Information retrieval is a field concerned with the structure, analysis, organization, storage, searching, and retrieval of information.*
- This class is about computational methods for the structure, analysis, organization, storage, searching, and retrieval of information.
  - And primarily about *text documents*.

# What is a Document?

- Examples:
  - web pages, email, books, news stories, scholarly papers, text messages, Word$^{TM}$, Powerpoint$^{TM}$, PDF, forum postings, patents, IM sessions, etc.

- Common properties:
  - Significant text content.
  - Some structure (e.g., title, author, date for papers; subject, sender, destination for email).

# Examples of Documents

# Documents vs. Database Records

- Database records are typically made up of well-defined fields (or *attributes)*
  - e.g. company names ,addresses ,account numbers, drug names, patent numbers, investigation file numbers.
- Easy to compare fields with well-defined semantics to queries in order to find matches.
- Our query has no fields and our documents have little structure.

# IR vs. Databases

**Information Retrieval**

- Data:
  - Semi-structured.
  - Heterogeneous.
  - Noisy.
- Unstructured or semi-structured queries.
- Natural language semantics.
- Infrequent off-line index changes.

**Databases**

- Data:
  - Structured.
  - Homogeneous.
  - Clean.
- Structured queries.
- Well-defined field semantics.
- Frequent on-line index changes.

# Generic Drugs – Illegal Activities by Manufacturers

# Comparing Text

- Determining whether a document matches a query is a fundamental problem of IR.

- Exact match is not enough:
  - Many different ways to state the same information
  - Documents may be relevant even when lacking some of the query terms.
  - Documents may be nonrelevant even if they contain all the query terms.

# Relevance

- What does it mean for a document to be *relevant*?
  - Simple definition: A relevant document contains information that a person was looking for when they submitted a query to the search engine.
  - Many factors influence a person's decision about what is relevant: e.g., task, context, novelty, style.
  - *Topical relevance* (same topic) vs. *user relevance* (everything else).

- How can we build an engine that retrieves relevant documents?

# Retrieval

- *Retrieval* models define a view of relevance.
- *Ranking algorithms* used in search engines are based on retrieval models.
- Most models describe statistical properties of text rather than linguistic properties.
  - i.e. counting simple text features such as words.
  - Statistical approach started with Luhn in the '50s.
  - Linguistic features can be part of a statistical model.

# Evaluation

- How do we know whether the engine is doing a good job of finding relevant documents?
  - *Evaluation* is experimental procedures and measures for comparing system output with user expectations.
  - IR evaluation methods now used in many fields.
  - *Recall* and *precision* are examples of effectiveness measures.

# Not Just Documents

- New applications increasingly involve media.
  - e.g. video, photos, music, speech
- Like text, contents is difficult to describe and compare.
  - text may be used to represent them (e.g. tags).
- IR approaches to search and evaluation are appropriate.

# Dimensions of IR

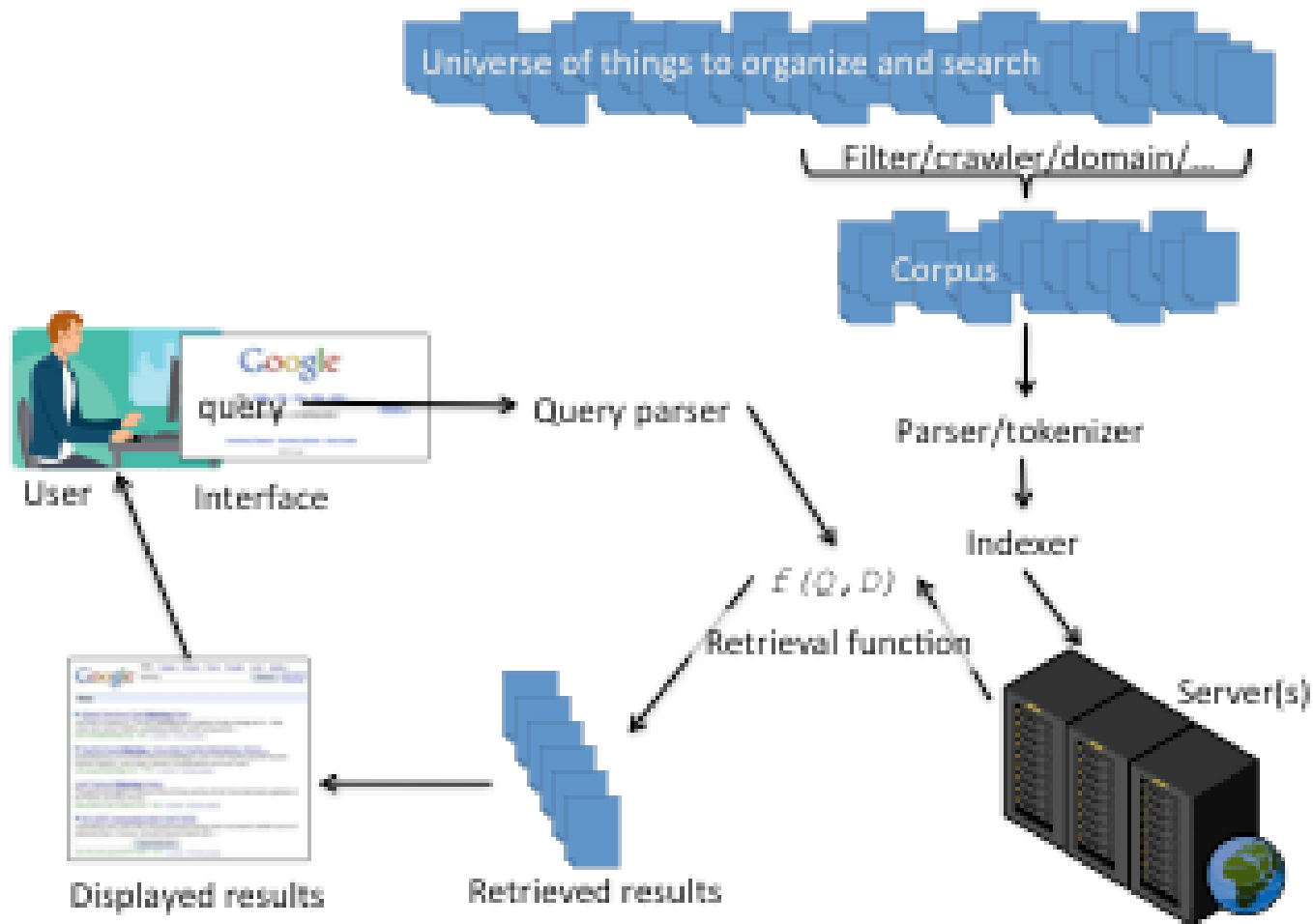| Content | Applications | Tasks |
|---------|-------------|-------|
| Text | Web search | Ad hoc search |
| Images | Vertical search | Filtering |
| Video | Enterprise search | Classification |
| Scanned docs | Desktop search | Question answering |
| Audio | Forum search | |
| Music | P2P search | |
| | Literature search | |

# IR Tasks

- Ad-hoc search:
  - Find relevant documents for an arbitrary text query.
- Filtering:
  - Identify relevant user profiles for a new document.
- Classification:
  - Identify relevant labels for documents.
- Question answering:
  - Give specific answer to a question.

# IR and Search Engines

- A search engine is the practical application of information retrieval techniques to large scale text collections.

- Relevance, retrieval, evaluation are issues.

- So are users and information needs, performance, coverage, updating, scalability, adaptability, and ability to handle specific problems (like spam).

# Components of a Search Engine

# Building a Search Engine

- Text processing and indexing.
  - Parsing; tokenizing; stopping and stemming; inverting indexes; scalability; index updates.
- Query processing and ranking.
  - Query languages; index look-up; retrieval models; features; relevance feedback; user interaction.
- Evaluation.
  - Effectiveness at performing task; querying speed; user satisfaction.

# Course Overview

- This course is about information retrieval in practice: the application of IR to search engine design and implementation.

- Course project:
  - Design and implement a small search engine capable of indexing and searching Wikipedia pages.
  - Evaluate its performance over provided queries.
  - Add something interesting to it.