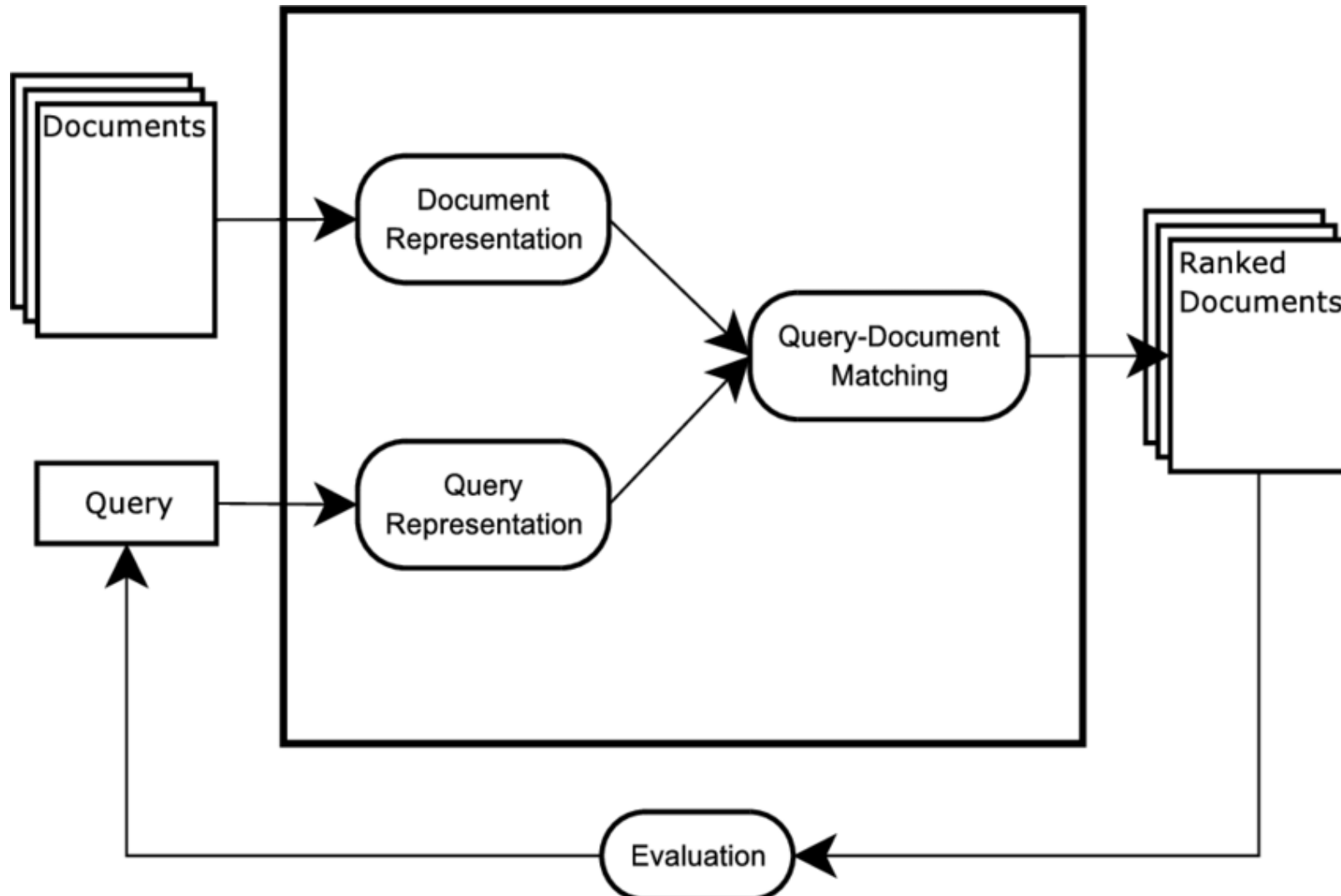# Information Retrieval

Amin Nazari

Spring 2025

# Outline

- What is Information Retrieval?
- IR applications
- IR vs Information Storage and Retrieval
- IR vs Data retrieval (DBMS)
- Why study it?
- Reference
- Outlines
- Grading

# What is 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).

- Unstructured data types
    - Text
    - Audio
    - Image
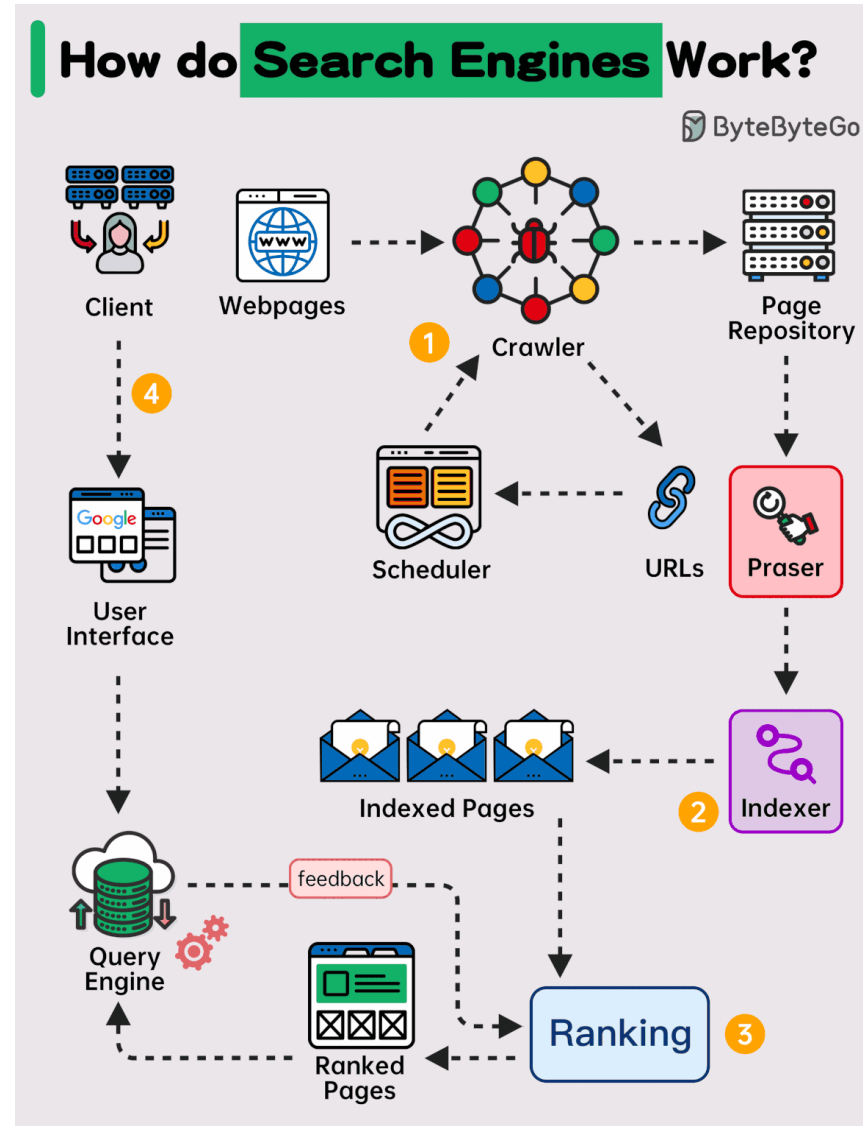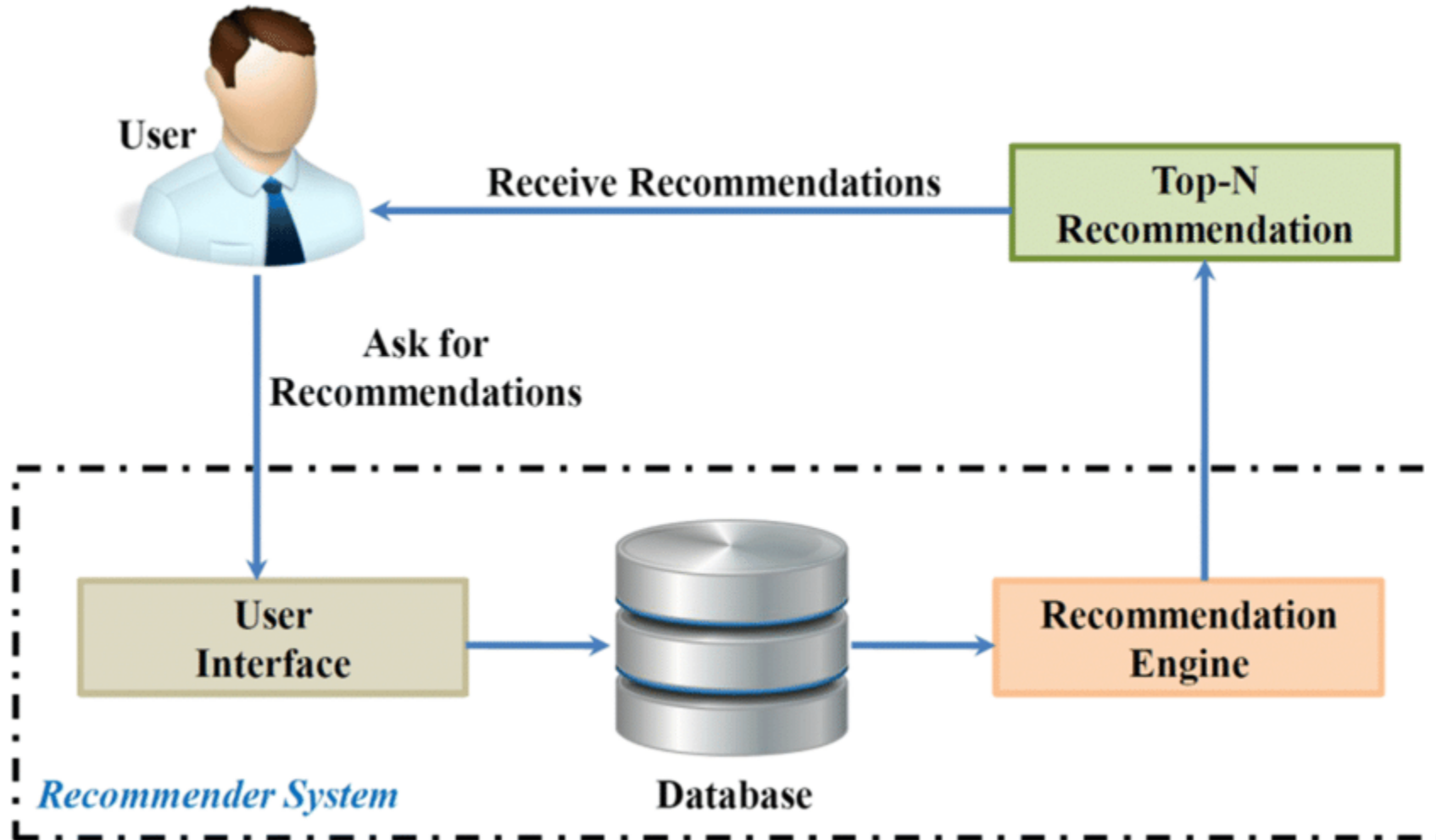    - Video
- Why text?

# IR schema

# Applications of IR

- Search engine
- Recommender systems (News, Movies, Posts, Usres, Books, Tour, Music, Drug, …)
- Documentation management (Digital Library)
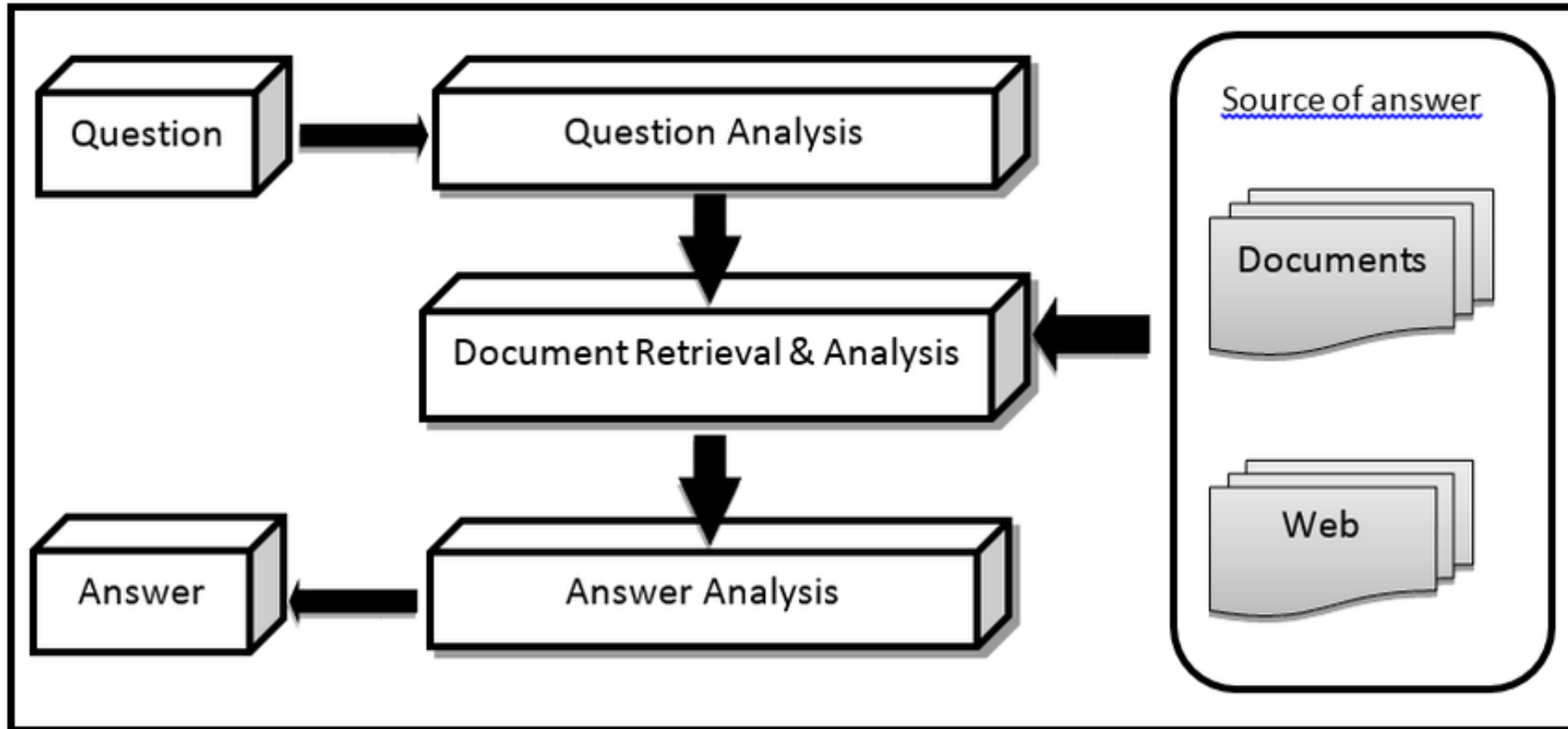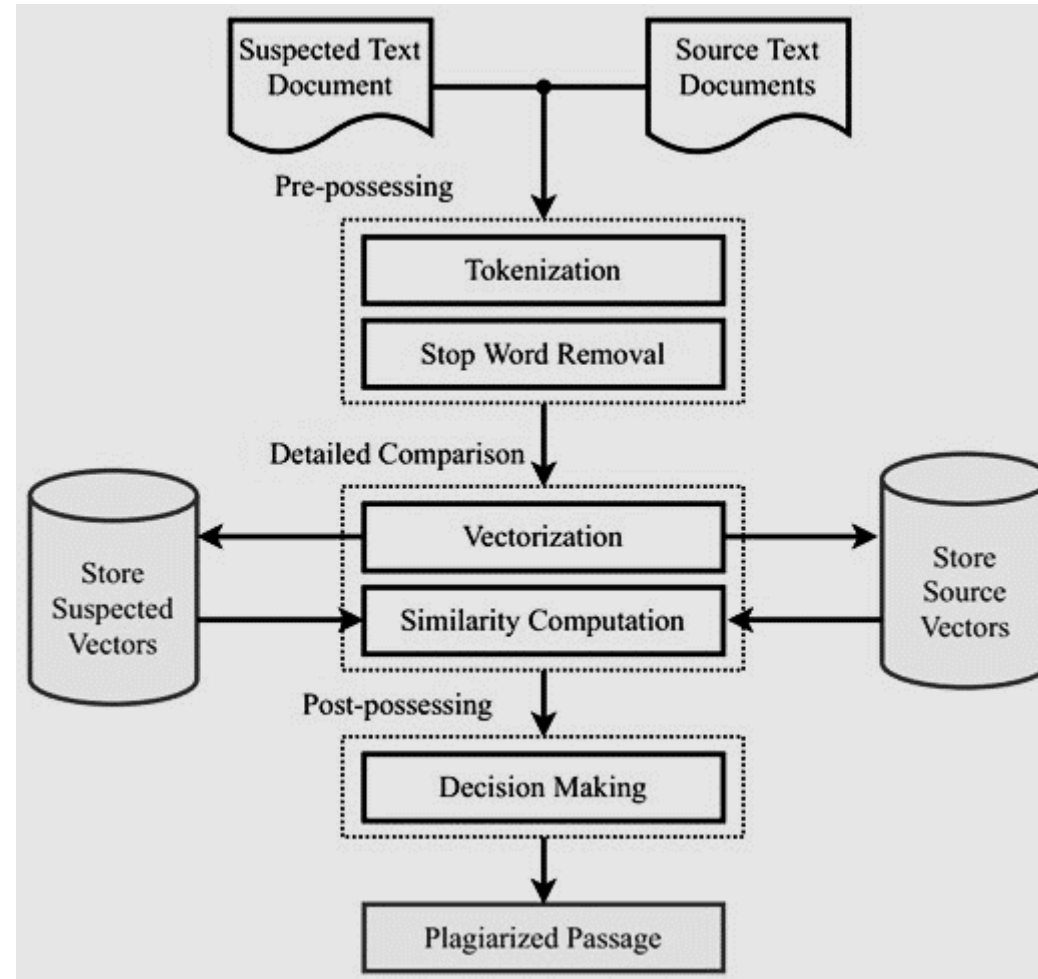- Question Answering Systems
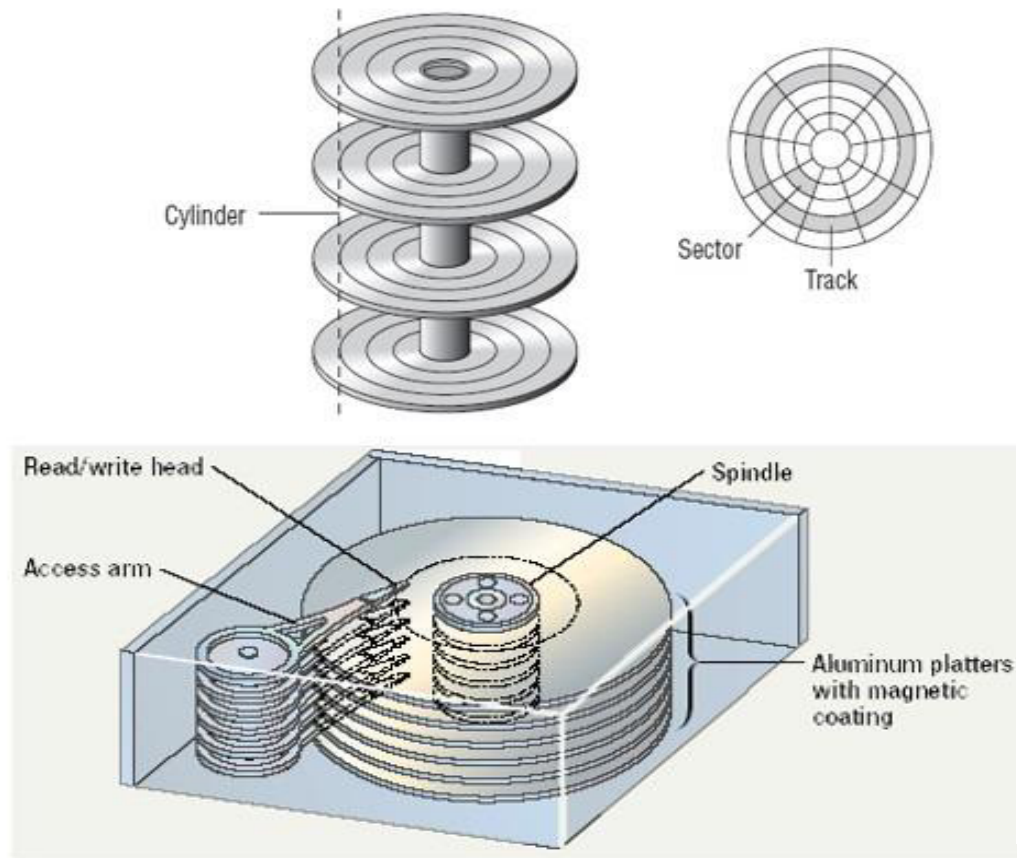- Plagiarism Detection

# Search Engines



How do Search Engines Work?

# RecSys

# Question Answering Systems

# Plagiarism Detection

# Information Storage and Retrieval

# Information Storage and Retrieval



label

polycarbonate plastic base

protective layer

adhesive layer

metal reflective layer

phase change layer

polycarbonate plastic base

spiral groove molded into the plastic base layer

laser    detector

a spot of laser altered phase change layer that represents the digital code stored on the DVD

**DVD**

**HD  DVD**

■Single-sided, single-layer

4. 7GB

15 GB

0. 6mm

0. 6mm

Bonding layer

Bonding layer

■Single-sided, dual-layer

8. 5GB

30 GB

0. 6mm

0. 6mm

switch   L1

switch   L1

L0

L0

Space  layer :55μm

Space  layer :20μm

# Information Storage and Retrieval

# Files disadvantages

- Distinguished and Isolated Data

- Data Duplication / Data Redundancy

- Data Protection

- Issues with Transactions –  ACID (Atomicity, Consistency, Isolation, and Durability)

- Concurrent issues

# DR vs IR

| Feature | Database System | Information Retrieval System |
|---|---|---|
| Data Type | Structured | Unstructured/Semi-structured |
| Query Language | SQL | Free-text/Natural Language |
| Result Matching | Exact | Relevance-based |
| Indexing | Structured (e.g., B-trees) | Inverted Index |
| Result Precision | Deterministic | Probabilistic |
| Examples | MySQL, PostgreSQL | Google, Elasticsearch |

# Why study IR?

- The most important problems in the domain of natural language processing (NLP)

- Hot topic research

- IR in LLM

- LLM in IR

- The roll of RecSys in e-commerce

# Outlines

| Theory | Practical |
|---|---|
| Text Preprocessing<br>Boolean and vector-space retrieval models<br>Evaluation and interface issues<br>Document clustering and classification<br>Traditional and machine learning-based ranking approaches<br>Recommender System<br>Web scrappy | Python<br>Numpy<br>Pandas<br>MatplotLib<br>Nltk<br>Regex (Regular Expressions)<br>NLTK (Natural Language Toolkit)<br>TextBlob |

# Intro course

- Introduction to Information Retrieval (https://nlp.stanford.edu/IR-book/newslides.html)
- CS 276: Information Retrieval and Web Search (stanford.edu)
- Grading:
  - Assignments: 25
  - Midterm: 25
  - Presentation: 10
  - Final: 40