



CSE 512: Distributed and Parallel Data Systems

Lecture 18

Instructor: Mohamed Sarwat

Teaching Evaluation

<https://fultonapps.asu.edu/eval/>



Google



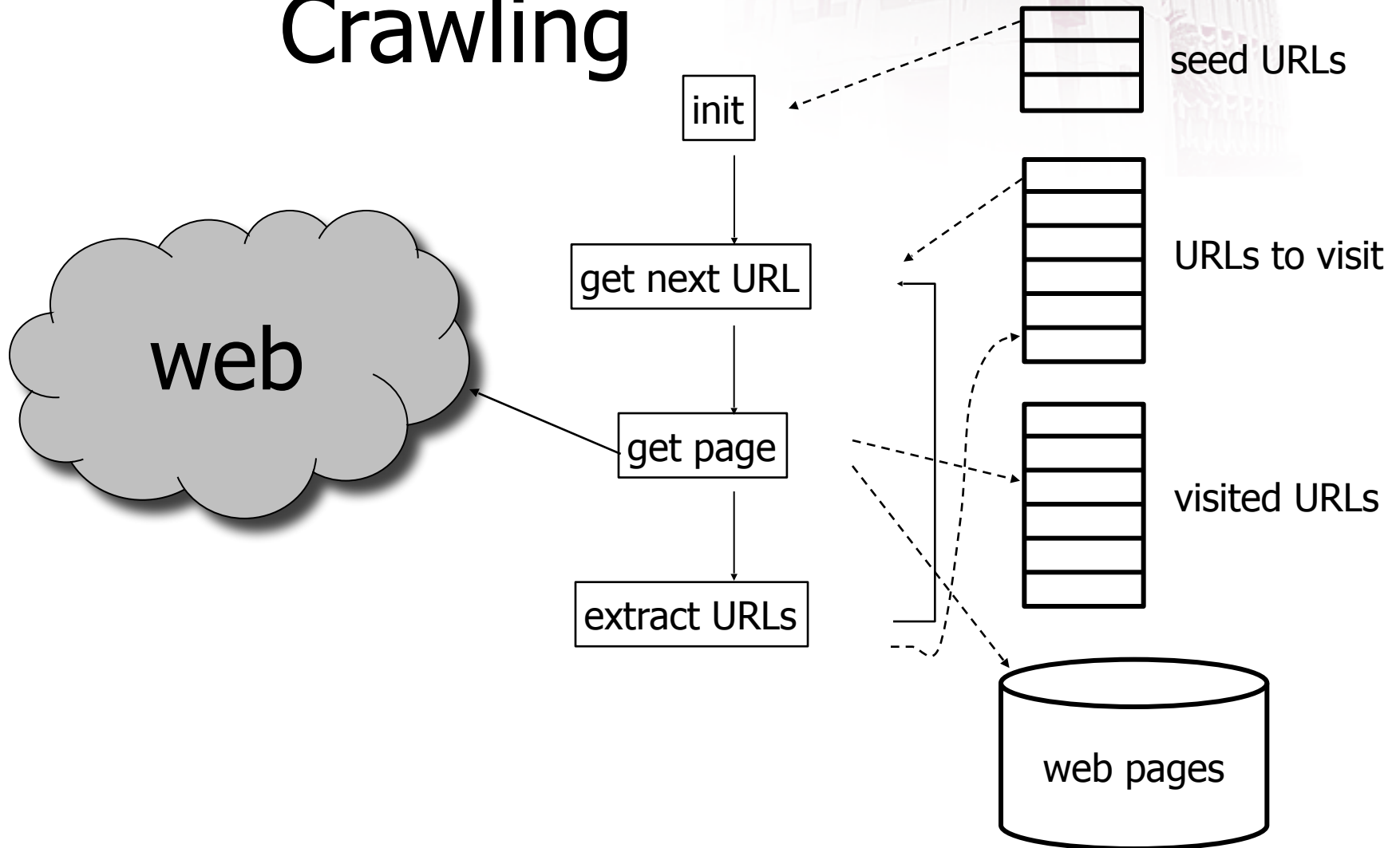
Google Search

I'm Feeling Lucky

Web Search Engine

- Crawling
- Indexing
- Computing ranking features
- Serving queries

Crawling



Issues

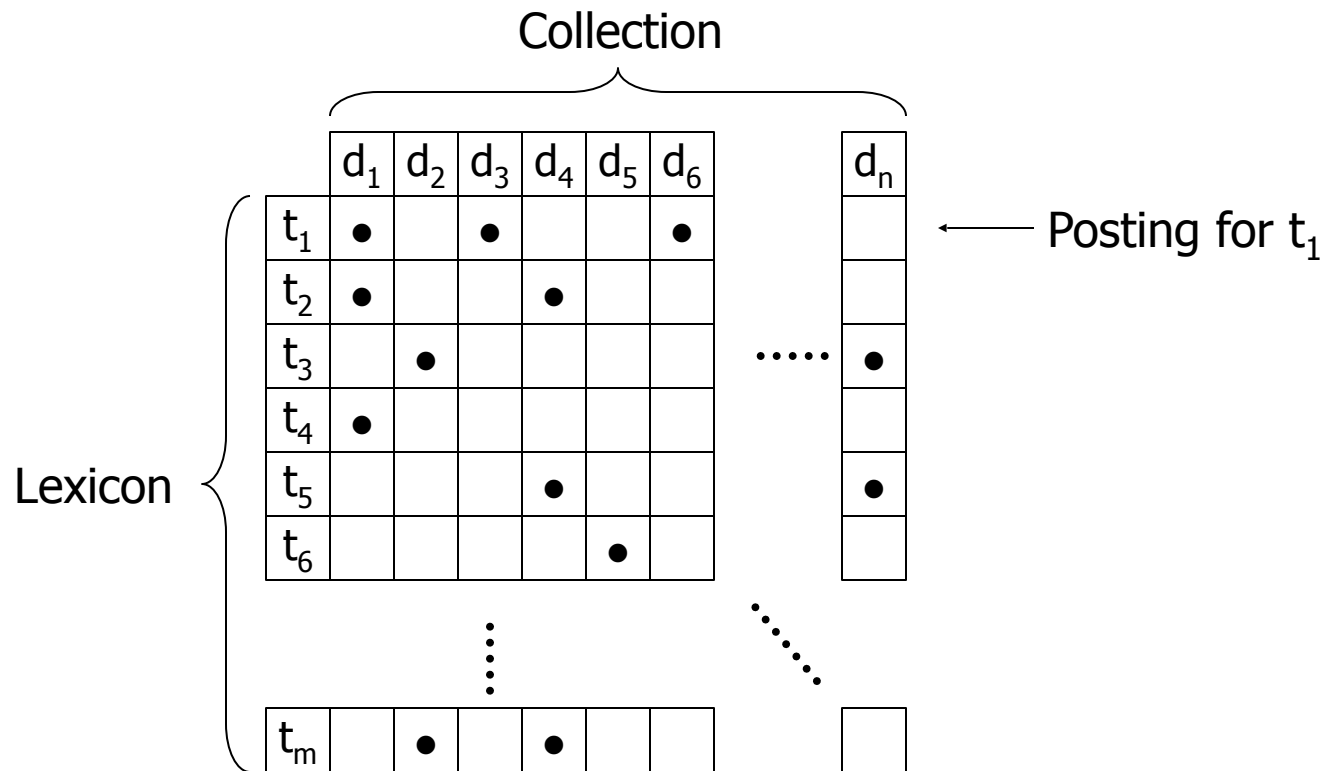
- Scope and freshness
 - Not enough space/time to crawl “all” pages
 - Page importance, quality, and update frequency
 - Site mirrors and (near) duplicate pages
 - Dynamic content and crawler traps
- Load at visited web sites
 - Rules in robots.txt
 - Limit number of visits per day
 - Limit depth of crawl

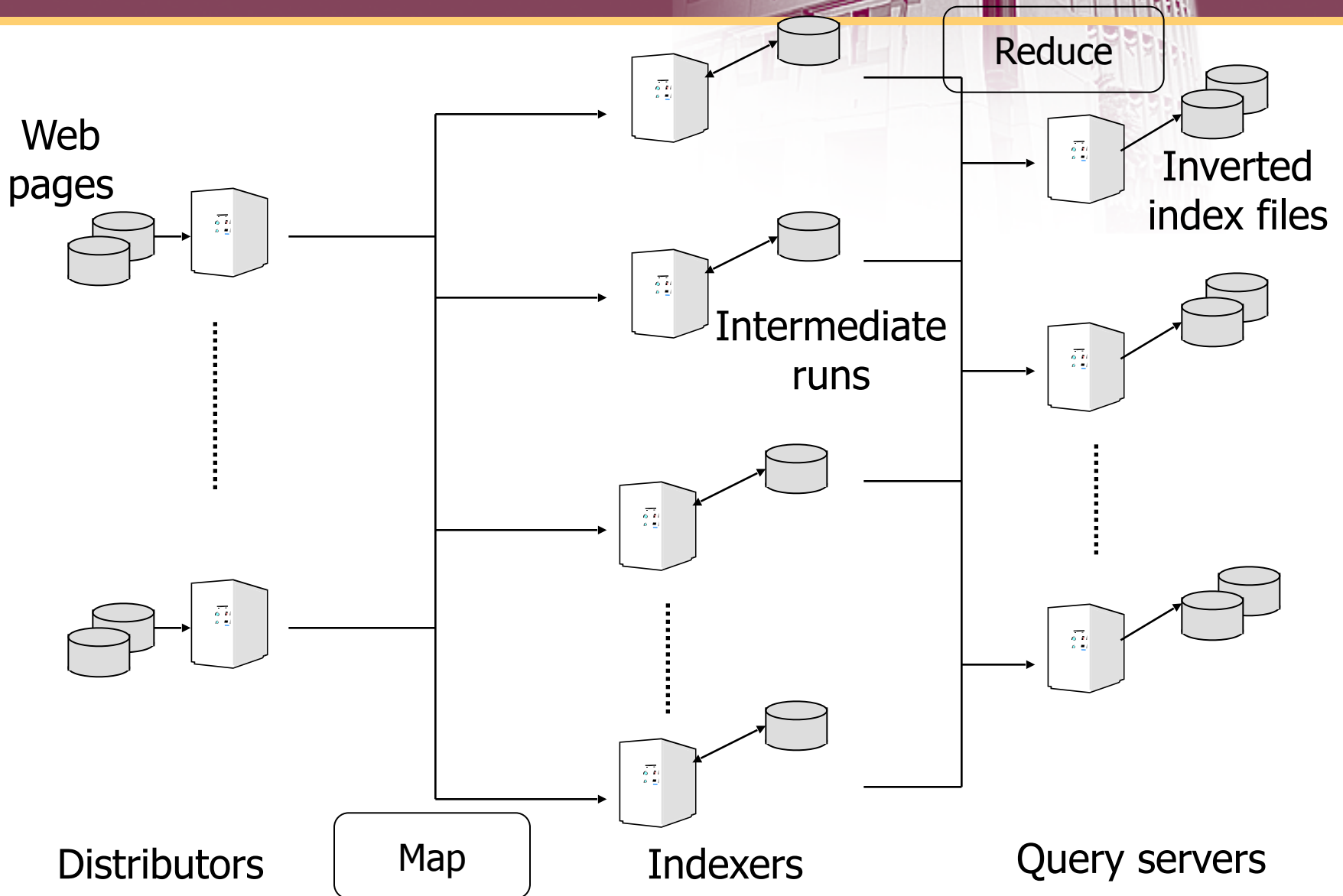
Issues

- Load at crawler
 - Variance of fetch latency/bandwidth
 - Parallelization and scalability
 - Multiple agents
 - Partitioning URL lists
 - Communication between agents
 - Recovering from agent failure

Indexing

- Build term-document index

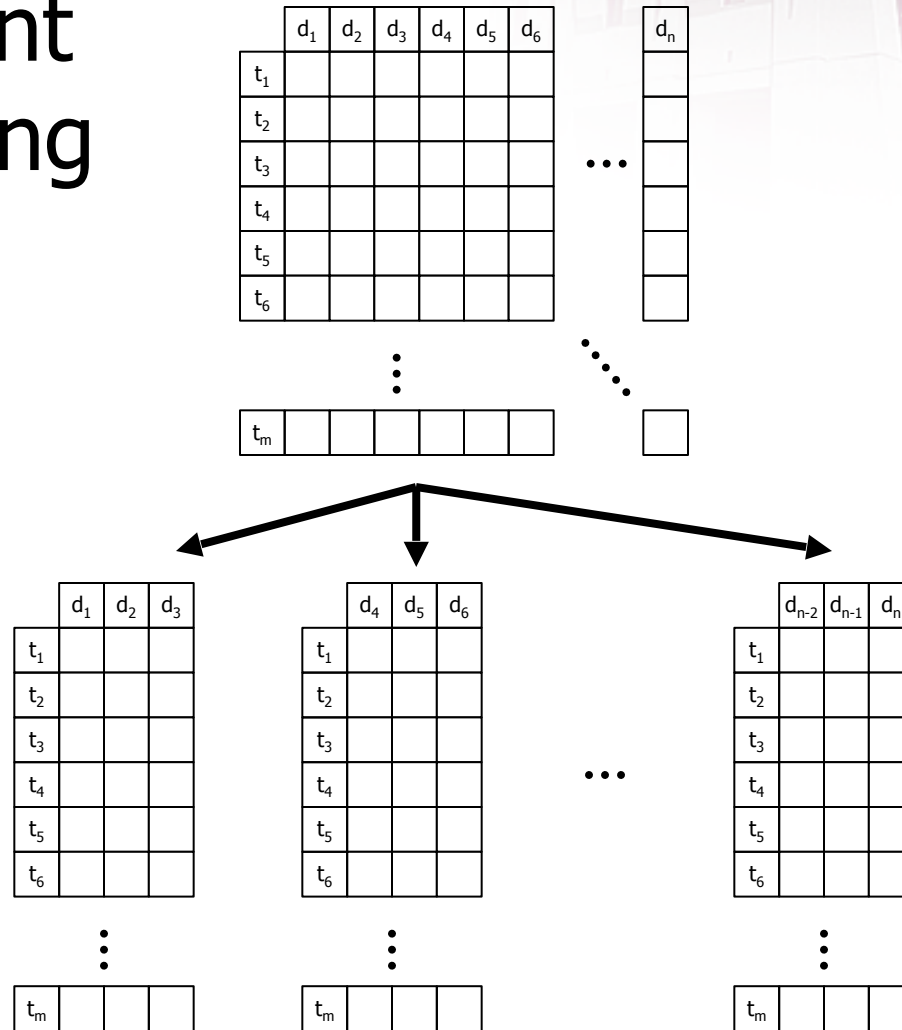




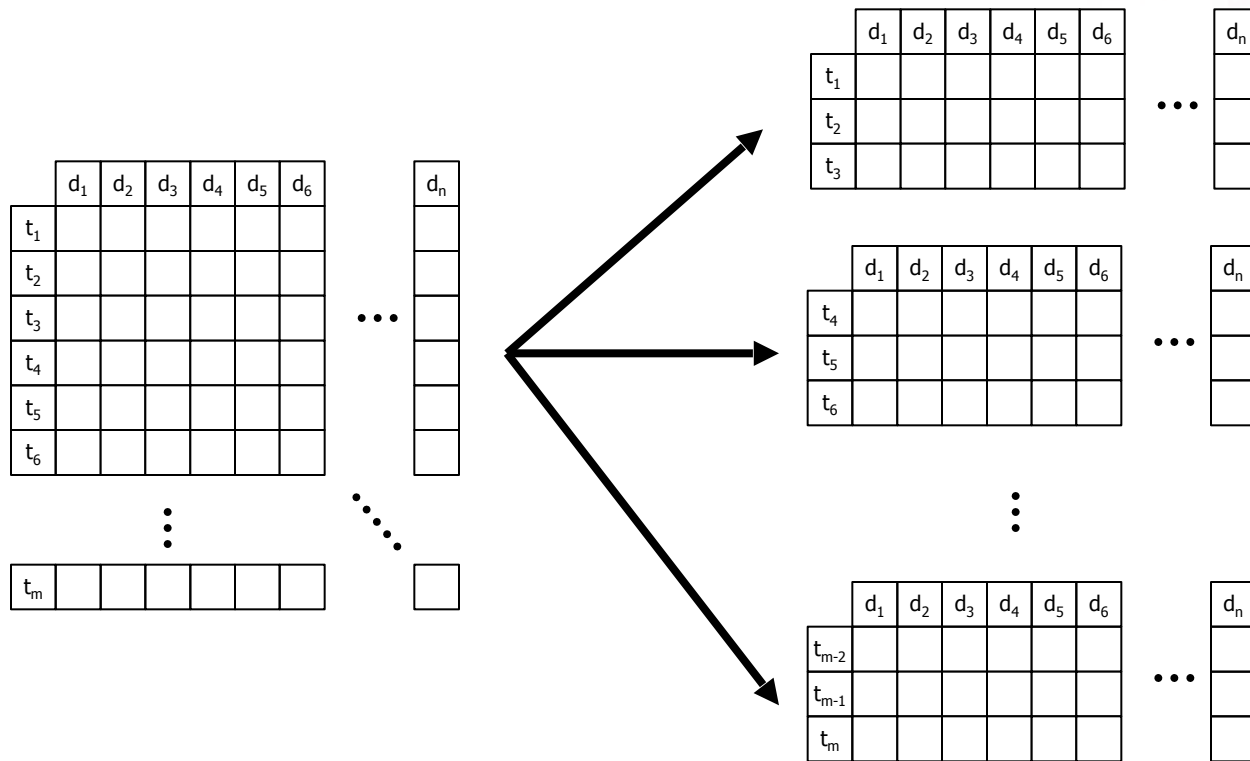
Issues

- Index partitioning
 - Efficient query processing
 - Query routing
 - Result retrieval

Document Partitioning



Term Partitioning



Document Partitioning

- Split the collection of documents
- Advantages
 - Easy to add new documents
 - Load balanced
 - High processing throughput
- Disadvantages
 - Communication with all query servers

Term Partitioning

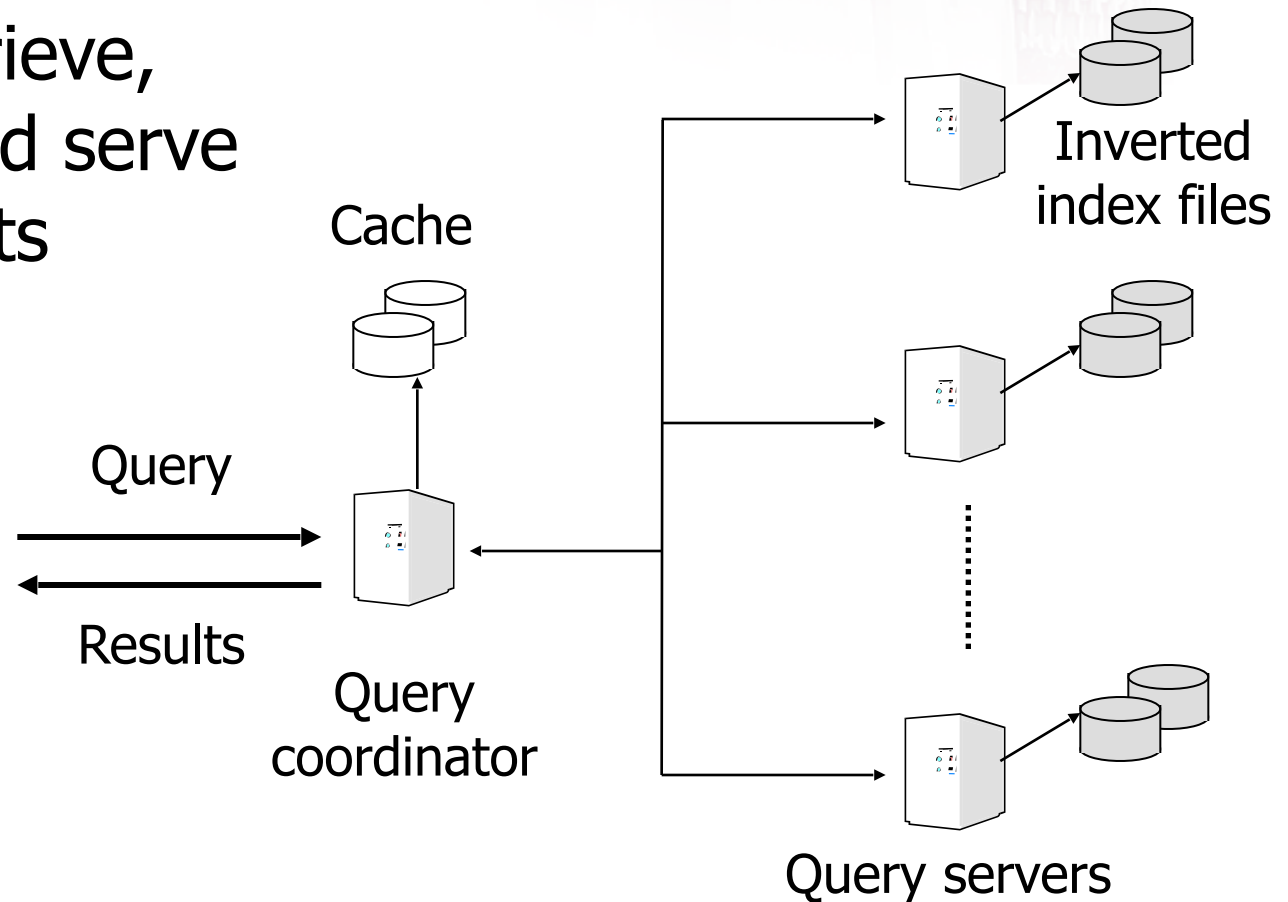
- Split the lexicon
- Advantages
 - Reduced communication with query servers
- Disadvantages
 - More processing before partitioning
 - Adding new documents is hard
 - Load balancing is hard
 - Processing throughput limited by query length

Ranking Feature Computation

- Parallel/distributed computation tasks
 - Text/language processing
 - Document classification/clustering
 - Web graph analysis

Query Processing

- Locate, retrieve, process, and serve query results



Architecture

- Multiple sites connected by WAN
 - Site = coordinator + servers + cache
- Partitioning
 - Parallel processing
 - Distributed storage of data
 - E.g., index partitioning
- Replication
 - Availability
 - Throughput
 - Response time





Wrap-Up

What we Covered

- Database Systems
- Distributed Systems
- Data Fragmentation (Partitioning)
- Distributed Query Processing and Optimization
- Distributed Transaction Processing



What we Covered

- Replicated Data Management
- Reliable Data Management
- Parallel Database Systems
- Data Management in MapReduce Systems
- NoSQL Data Management System

Take-Away Message

- Distributed and Parallel Data Systems are very hot these days
 - Major Internet companies require such skills
- I have tried to give you the abstract concept and also expose you to new technologies
- No one size (database system) fits all