

UNIVERSITY of DELAWARE

Web Search Basics

(Search and Data Mining)

Hui Fang
Department of Electrical and Computer Engineering
University of Delaware

UNIVERSITY of DELAWARE

Characteristics of Web Data

- The size of Web data is huge.**
 - Surface: static pages
 - Deep: dynamically generated pages
- Frequent updates**
- Variances in quality**
 - Web Spams
- Graph-based data**

UNIVERSITY of DELAWARE

General Challenges in Web Search

- Scalability
 - How to handle the size of the Web?
- Addressing the dynamics of the Web
 - Some pages may be updated very quickly
 - New pages are constantly created
- Dealing with low quality information
- Utilizing the additional information (such as links) to improve search quality

Parallel query processing (MapReduce)

Crawling/
Redundancy check

Spam detection

Link analysis / multi-feature ranking

UNIVERSITY of DELAWARE

Component I: Crawler

- Basic steps of crawling
 - Start with a set of seed pages
 - Fetch pages from the Web
 - Parse the fetched pages and add the hyperlink to the list to be crawled
 - Repeat the previous steps until the list to be crawled is empty.

4

UNIVERSITY of DELAWARE

More Details about Crawling

- Breadth-First
- Parallel crawling
- Focused crawling
- Incremental/repeated crawling

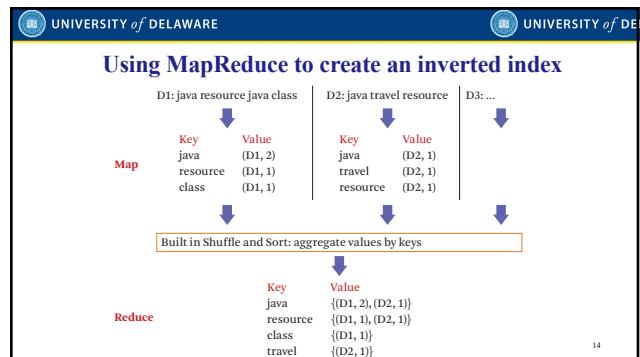
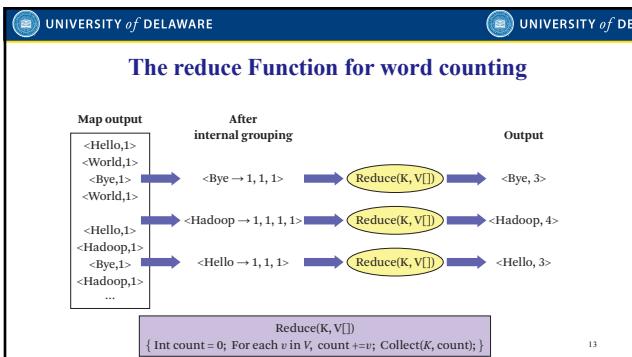
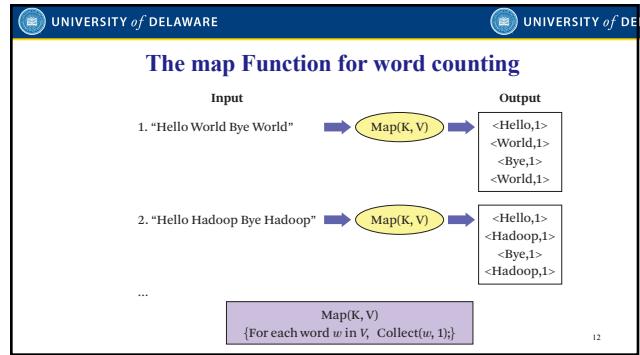
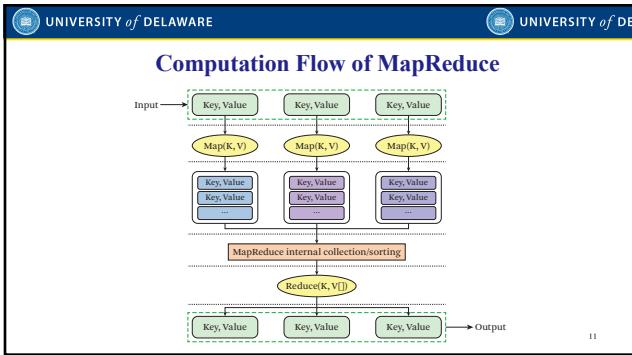
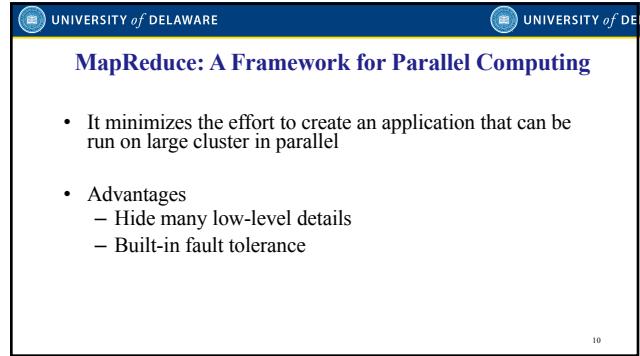
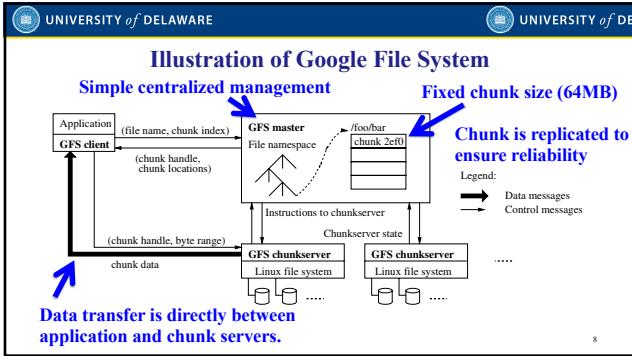
5

UNIVERSITY of DELAWARE

Component II: Indexing

- Standard IR techniques are the basis, but insufficient
 - Scalability
 - Efficiency
- Google's innovations on Web-scale data processing
 - Google File System
 - Storing the indexes on multiple machines
 - MapReduce
 - Creating the indexes in parallel
 - BigTable

7



UNIVERSITY of DELAWARE

Pseudocode for Inverted Index Construction

```

function MAP(docid n, doc d)
    H ← new ASSOCIATIVEARRAY
    for term t ∈ d do
        H[t] ← H[t] + 1
    end for
    for t ∈ H do
        EMIT(t, [n, H[t]])
    end for
end function

function REDUCE(term t, postings [(a1, f1), (a2, f2), ...])
    P ← new LIST
    for all (a, f) do
        APPEND(P, (a, f))
    end for
    SORT(P)
    EMIT(t, P)
end function

```

15

UNIVERSITY of DELAWARE

Component III: Retriever

- Traditional IR models are not enough
 - Information needs are different: navigational vs. informational
 - Documents have additional information
 - Information quality varies a lot
- Specific techniques for Web search
 - Exploiting links
 - Spelling correction
 - Spam detection
 - Leveraging more features such as clickthrough data
- In general, rely on machine learning to combine all kinds of features.

17