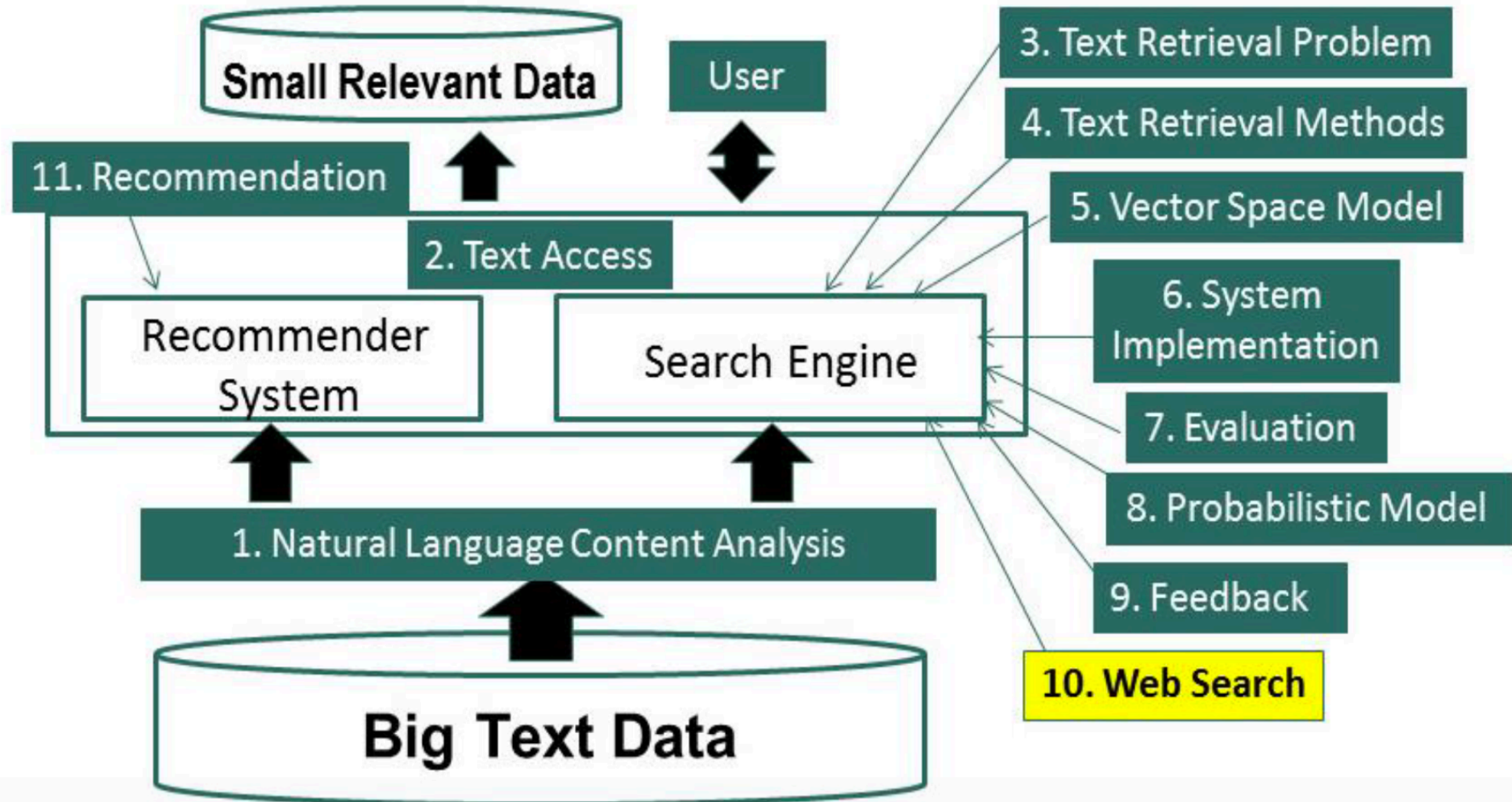# Information Retrieval & Text Mining
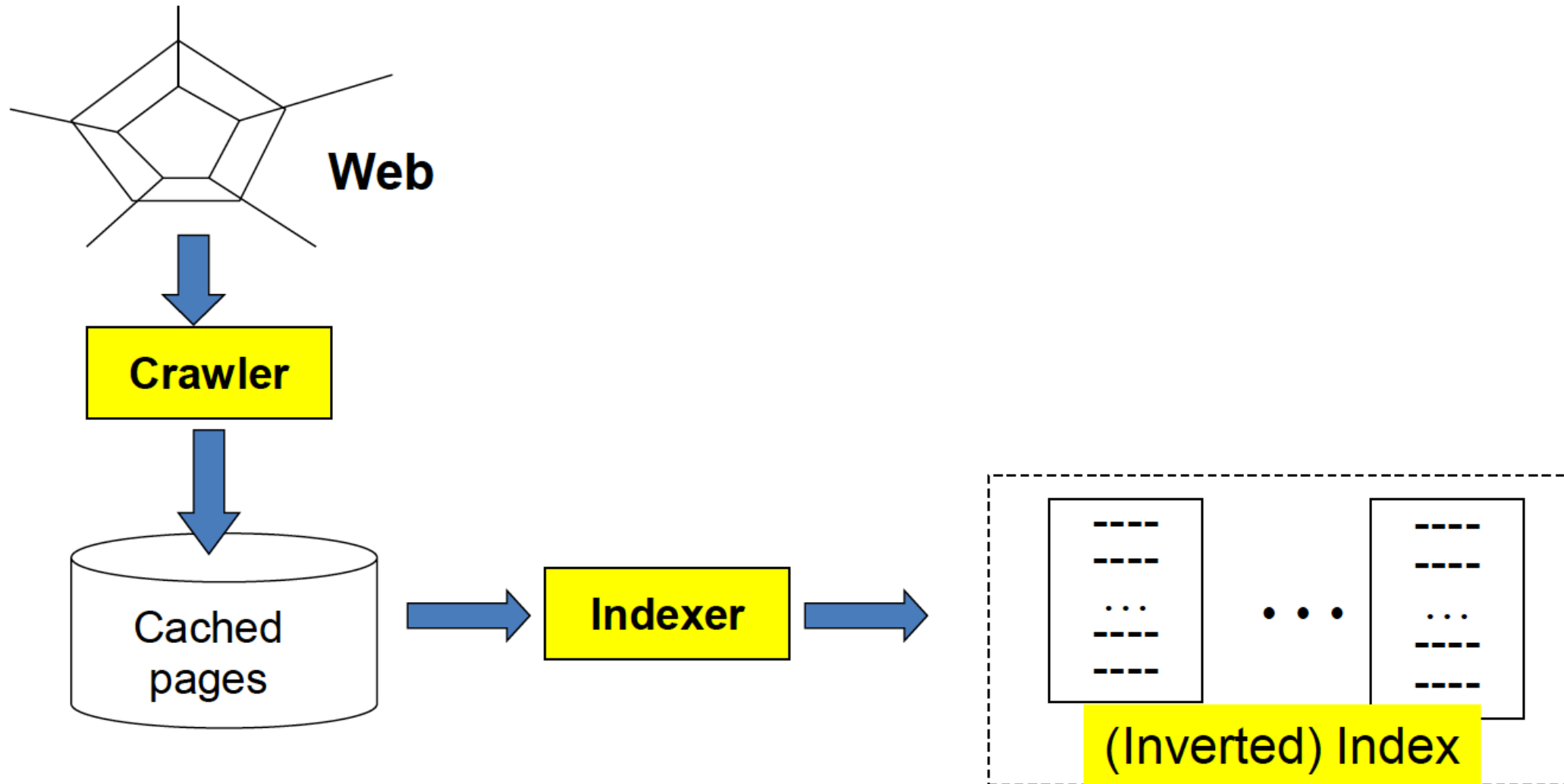
## Web Search
## Web Index

**Dr. Iqra Safder**
**Information Technology University**

# Course Schedule

# Basic Search Engine Technologies

# Overview of Web Indexing

- Standard IR techniques are the basis, but insufficient
  - Scalability
  - Efficiency
- Google's contributions:
  - Google File System (GFS): distributed file system
  - MapReduce: Software framework for parallel computation
  - Hadoop: Open source implementation of MapReduce

# GFS Architecture

Simple centralized management

Fixed chunk size (64 MB)

Chunk is replicated to ensure reliability

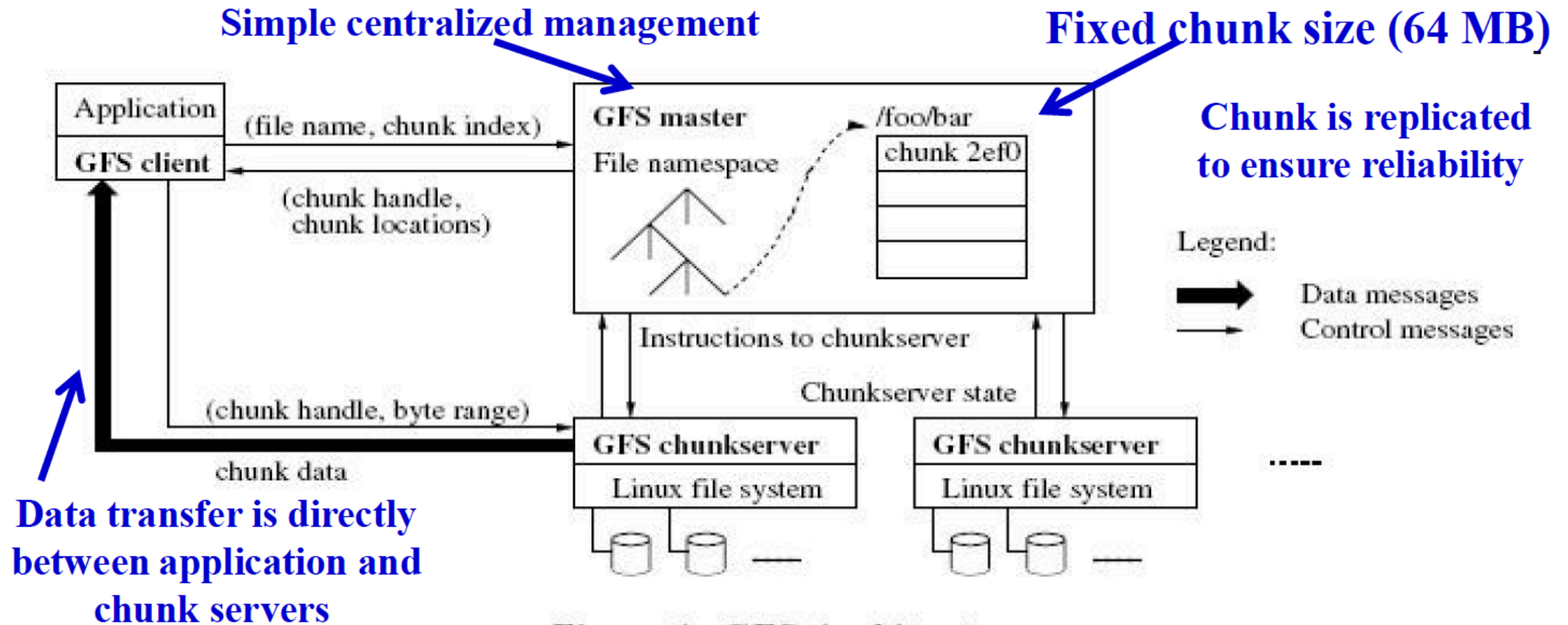Data transfer is directly between application and chunk servers

Figure 1: GFS Architecture

GHEMAWAT, S., GOBIOFF, H., AND LEUNG, S.-T. The google file system. In SOSP '03: Proceedings of the nineteenth ACM symposium on Operating systems principles (New York, NY, USA, 2003), ACM, pp. 29–43.
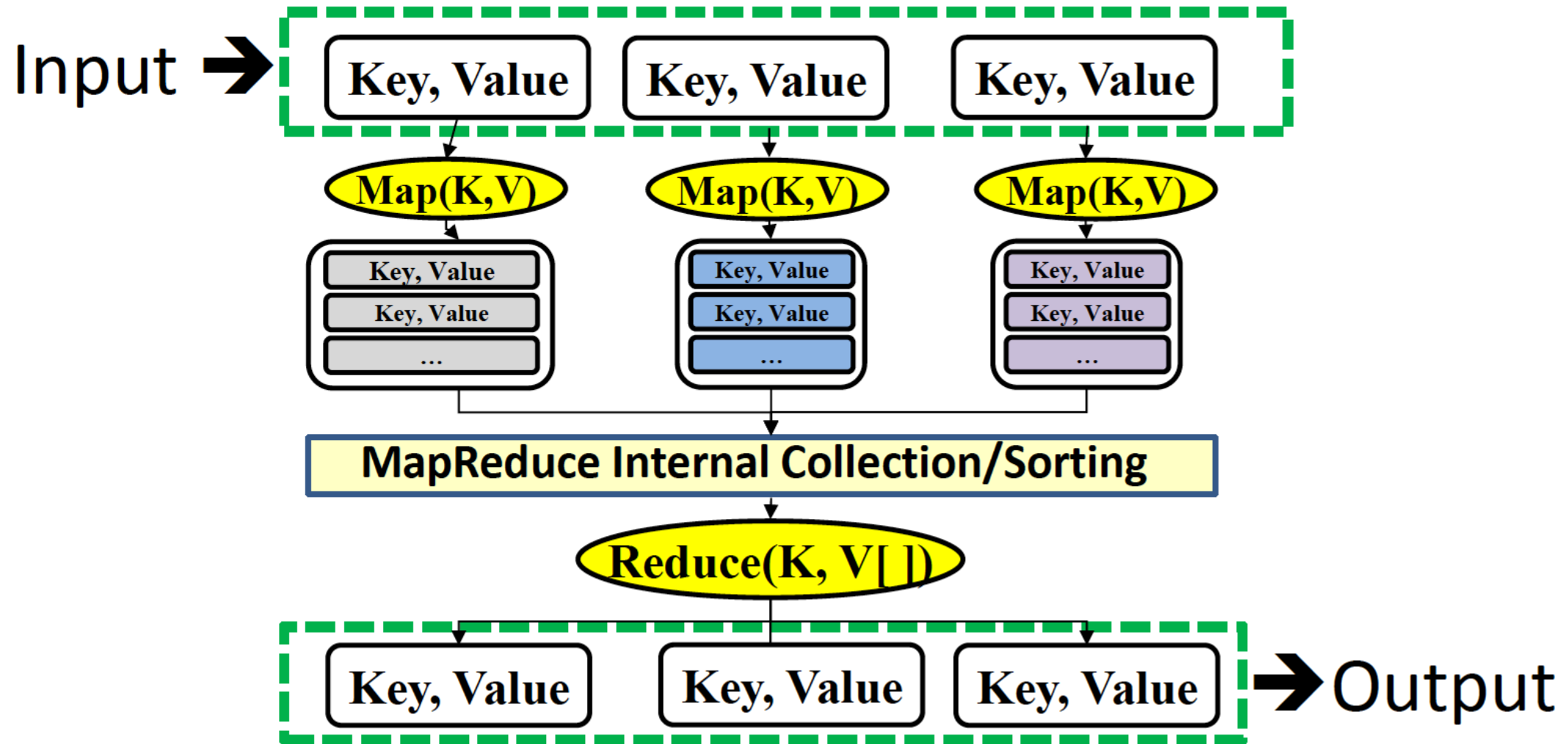
http://static.googleusercontent.com/media/research.google.com/en/us/archive/gfs-sosp2003.pdf

# MapReduce: A Framework for Parallel Programming

- Minimize effort of programmer for simple parallel processing tasks
- Features
  - Hide many low-level details (network, storage)
  - Built-in fault tolerance
  - Automatic load balancing

# MapReduce: Computation Pipeline

# Word Counting: Map Function

**Input**

1, "Hello World Bye World" ➡️ **Map(K,V)** ➡️
```
<Hello,1>
<World,1>
<Bye,1>
<World,1>
```

**Output**

2, "Hello Hadoop Bye Hadoop" ➡️ **Map(K,V)** ➡️
```
<Hello,1>
<Hadoop,1>
<Bye,1>
<Hadoop,1>
```

….

**Map(K, V)**
**{ For each word w in V,  Collect(w, 1);}**

# Word Counting: Reduce Function

**Map Output**

**After internal grouping**

**Output**

<Hello,1>
<World,1>
<Bye,1>
<World,1>

<Hello,1>
<Hadoop,1>
<Bye,1>
<Hadoop,1>
…

<Bye → 1, 1, 1>  →  **Reduce(K, V[ ])**  →  <Bye, 3>

<Hadoop → 1, 1, 1, 1>  →  **Reduce(K, V[ ])**  →  <Hadoop, 4>

<Hello → 1, 1, 1>  →  **Reduce(K, V[ ])**  →  <Hello, 3>
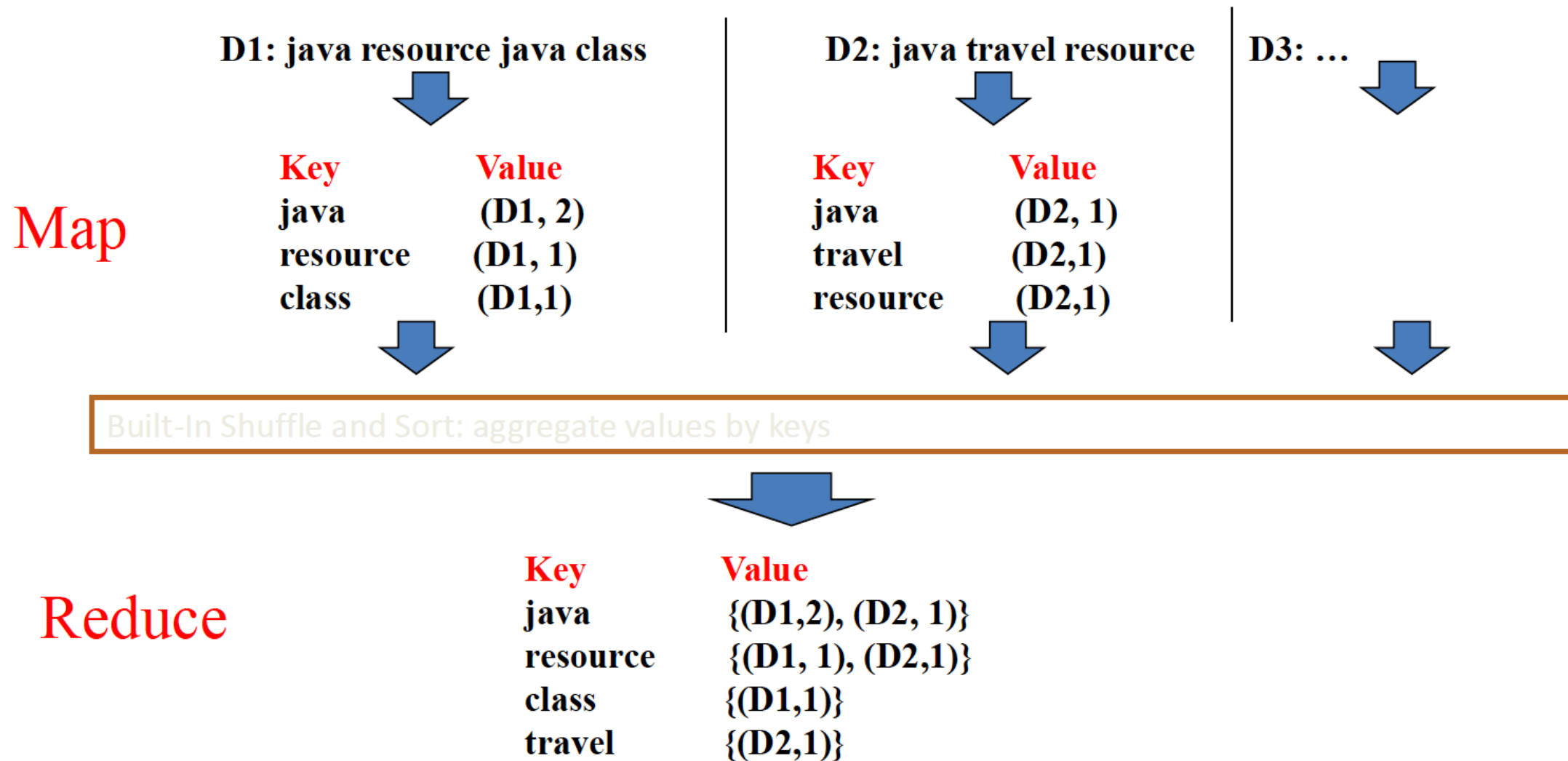
**Reduce(K, V[ ])**
**{ Int count = 0;  For each v in V,  count += v;   Collect(K, count); }**

# Inverted Indexing with MapReduce

**D1: java resource java class**                **D2: java travel resource**          **D3: …**

**Map**

| Key | Value |
|-----|-------|
| java | (D1, 2) |
| resource | (D1, 1) |
| class | (D1,1) |

| Key | Value |
|-----|-------|
| java | (D2, 1) |
| travel | (D2,1) |
| resource | (D2,1) |

Built-In Shuffle and Sort: aggregate values by keys

**Reduce**

| Key | Value |
|-----|-------|
| java | {(D1,2), (D2, 1)} |
| resource | {(D1, 1), (D2,1)} |
| class | {(D1,1)} |
| travel | {(D2,1)} |

…

Slide adapted from Jimmy Lin's presentation

# Inverted Indexing – Pseudo code

```
1: procedure Map(k, d)
2:     Initialize.AssociativeArray(H)
3:     for all t ∈ d do
4:         H{t} ← H{t} + 1
5:     for all t ∈ H do
6:         Emit(t, ⟨k, H{t}⟩)
1: procedure Reduce(t, [⟨k₁, f₁⟩, ⟨k₂, f₂⟩ ...])
2:     Initialize.List(P)
3:     for all ⟨k, f⟩ ∈ [⟨k₁, f₁⟩, ⟨k₂, f₂⟩ ...] do
4:         Append(P, ⟨k, f⟩)
5:     Sort(P)
6:     Emit(t, P)
```

# Summary

- Web scale indexing requires
  - Storing the index on multiple machines (GFS)
  - Creating the index in parallel (MapReduce)
- Both GFS and MapReduce are general infrastructures