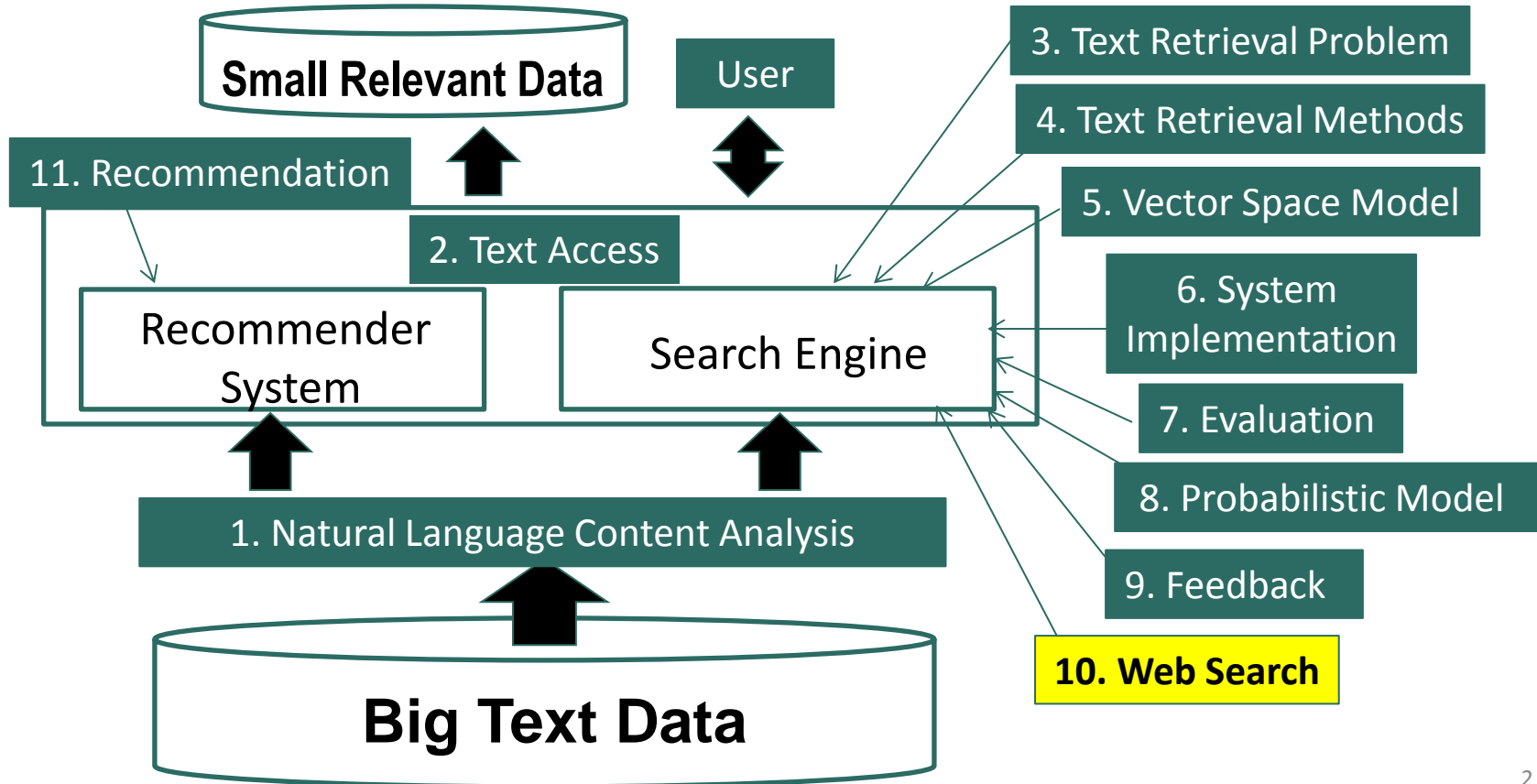# Text Retrieval and Search Engines
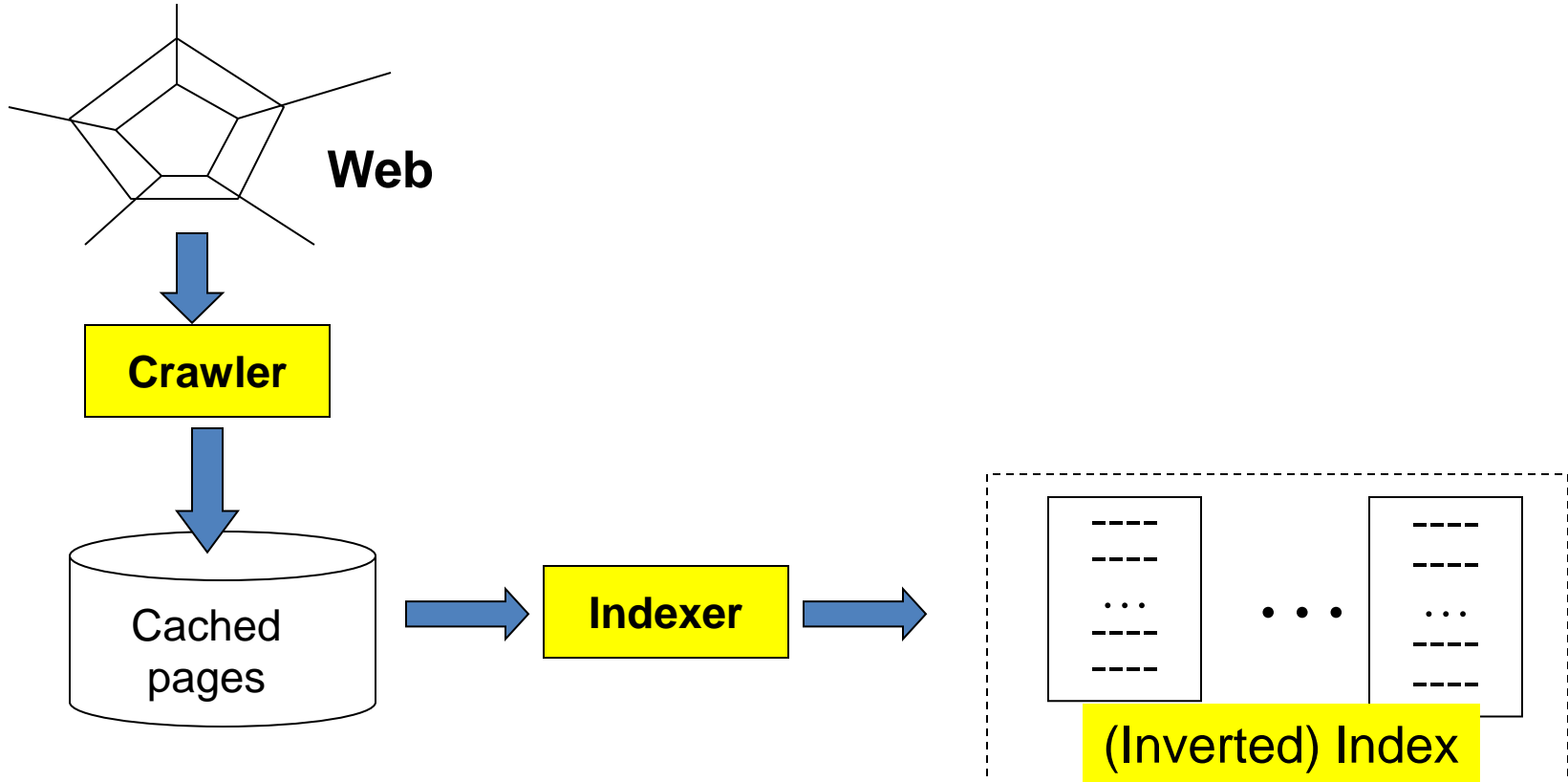
Web Search: Web Indexing

ChengXiang "Cheng" Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

# Web Search: Web Indexing

# Basic Search Engine Technologies

**Web**

**Crawler**

Cached pages

**Indexer**

- - - -
- - - -
- - -
· · ·
- - -
- - - -

- - - -
- - - -
- - -
· · ·
- - -
- - - -

(Inverted) Index

3

# 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**

| Application | (file name, chunk index) | GFS master |
|---|---|---|
| GFS client | | File namespace |

/foo/bar
chunk 2ef0

(chunk handle, chunk locations)

Legend:
→ Data messages
→ Control messages

Instructions to chunkserver

Chunkserver state

(chunk handle, byte range)

GFS chunkserver
Linux file system

GFS chunkserver
Linux file system

......

chunk data

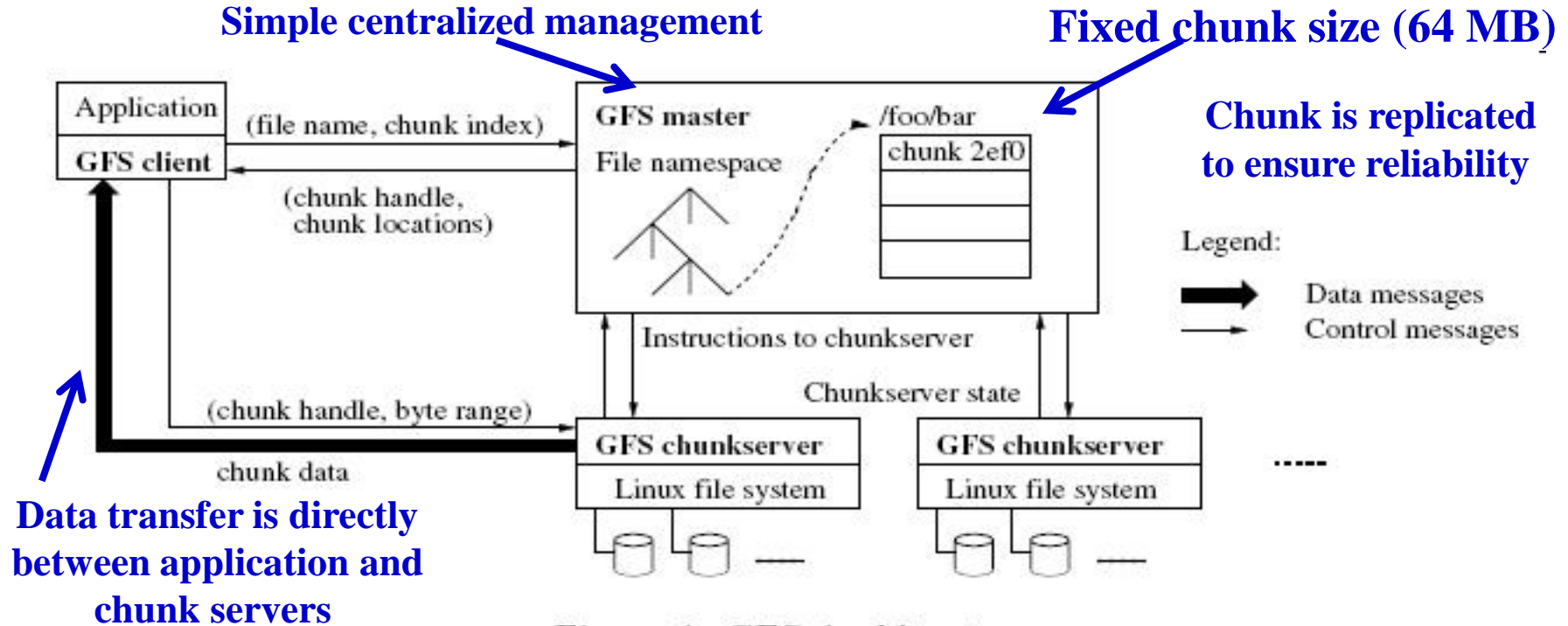**Data transfer is directly between application and chunk servers**
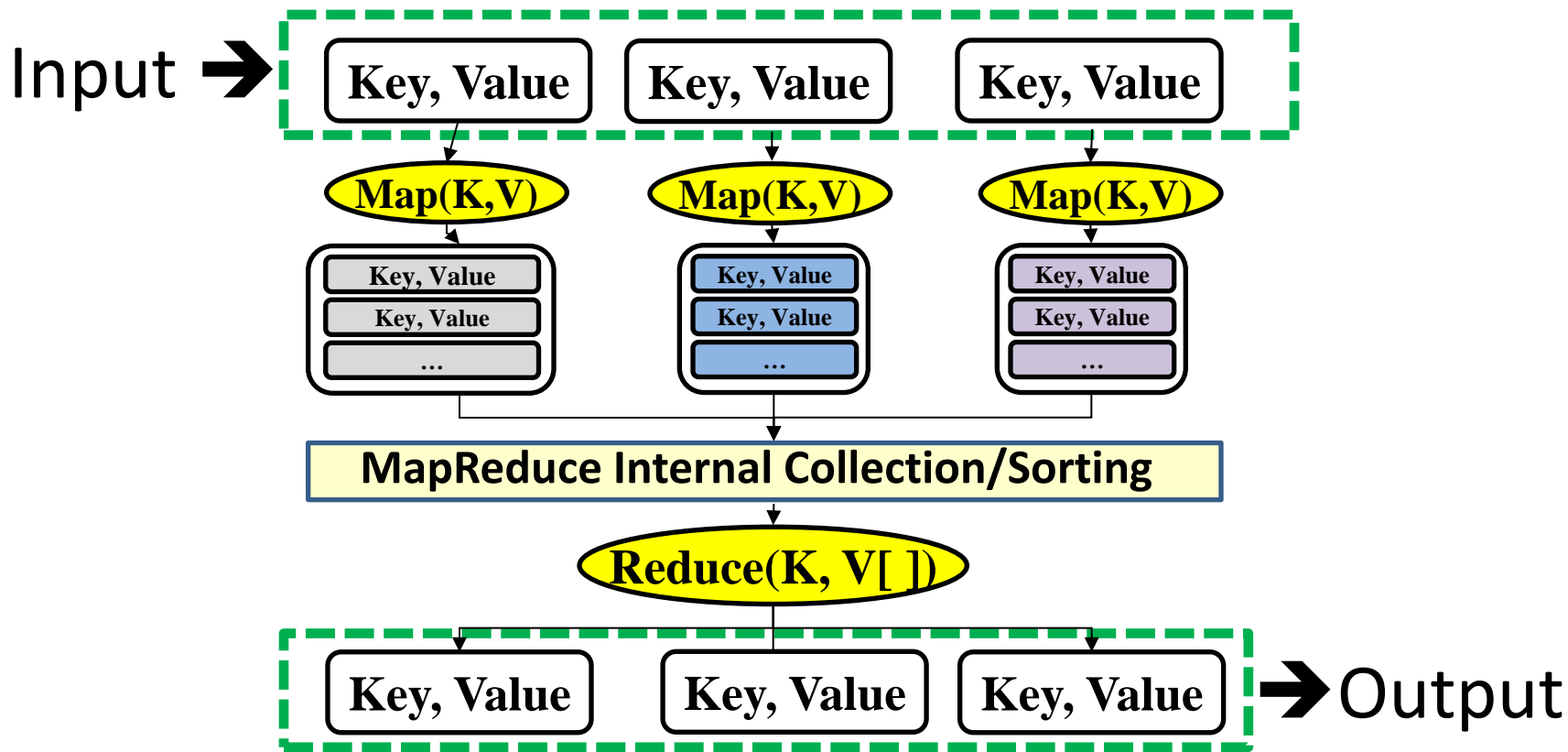
**Figure 1: GFS Architecture**

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

**Input: Text Data**

Hello World Bye World
Hello Hadoop Bye Hadoop
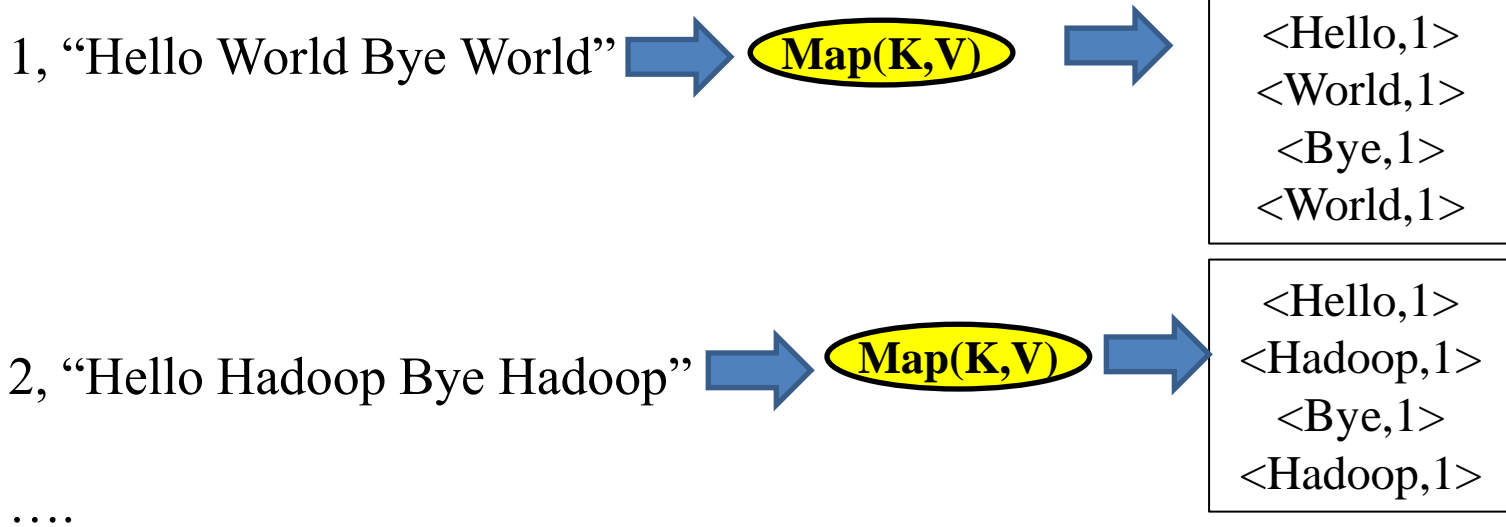Bye Hadoop Hello Hadoop
…  …

**Output:**
**Count of each word**

Bye 3
Hadoop 4
Hello 3
World 2
…

**How can we do this within the MapReduce framework?**

# Word Counting: Map Function

**Input**

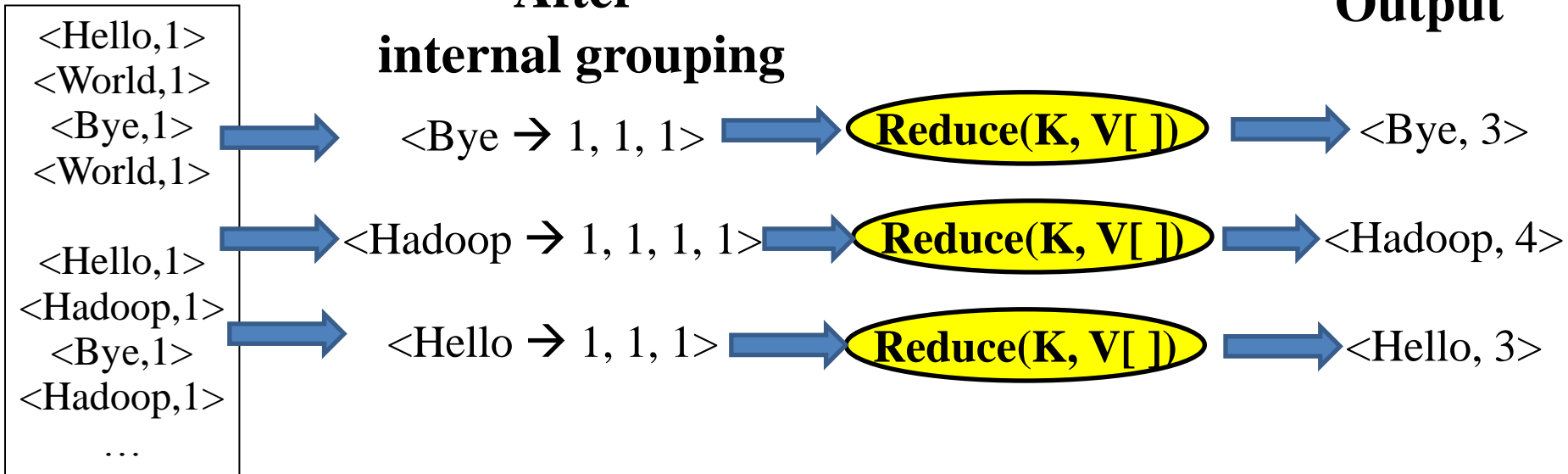**Output**

1, "Hello World Bye World" → **Map(K,V)** →

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

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**

| |
|---|
| <Hello,1> |
| <World,1> |
| <Bye,1> |
| <World,1> |
| |
| <Hello,1> |
| <Hadoop,1> |
| <Bye,1> |
| <Hadoop,1> |
| ... |

**After internal grouping**

<Bye → 1, 1, 1>

<Hadoop → 1, 1, 1, 1>

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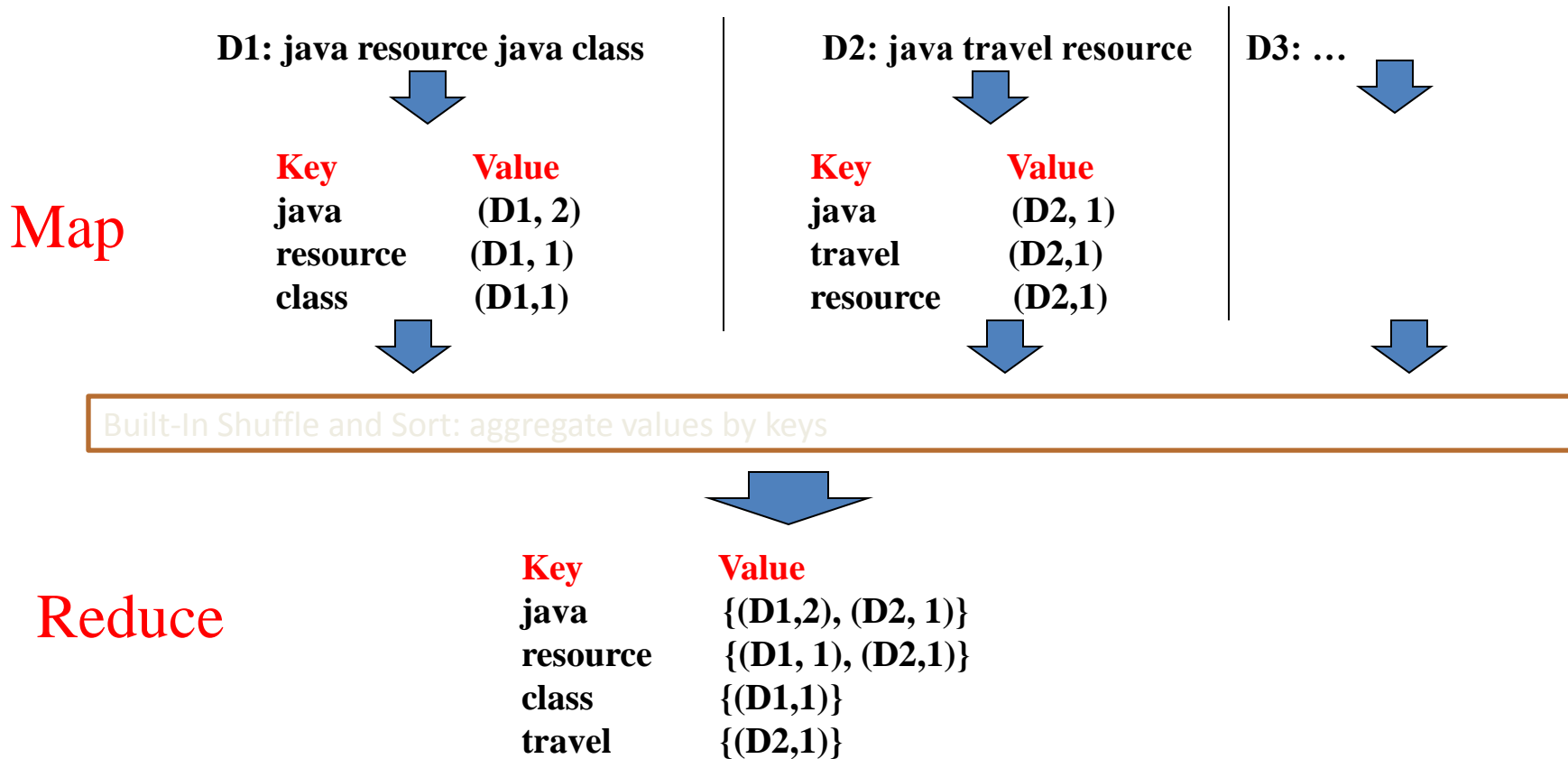
**Reduce(K, V[ ])**

**Output**

<Bye, 3>

<Hadoop, 4>

<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

11

# Inverted Indexing: Pseudo-Code

```
1: class Mapper
2:     procedure Map(docid n, doc d)
3:         H ← new AssociativeArray
4:         for all term t ∈ doc d do
5:             H{t} ← H{t} + 1
6:         for all term t ∈ H do
7:             Emit(term t, posting ⟨n, H{t}⟩)

1: class Reducer
2:     procedure Reduce(term t, postings [⟨a₁, f₁⟩, ⟨a₂, f₂⟩ . . .])
3:         P ← new List
4:         for all posting ⟨a, f⟩ ∈ postings [⟨a₁, f₁⟩, ⟨a₂, f₂⟩ . . .] do
5:             Append(P, ⟨a, f⟩)
6:         Sort(P)
7:         Emit(term t, postings 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