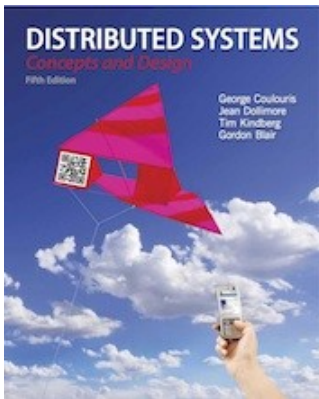


# Slides for Chapter 21: Designing Distributed Systems: Google Case Study

---

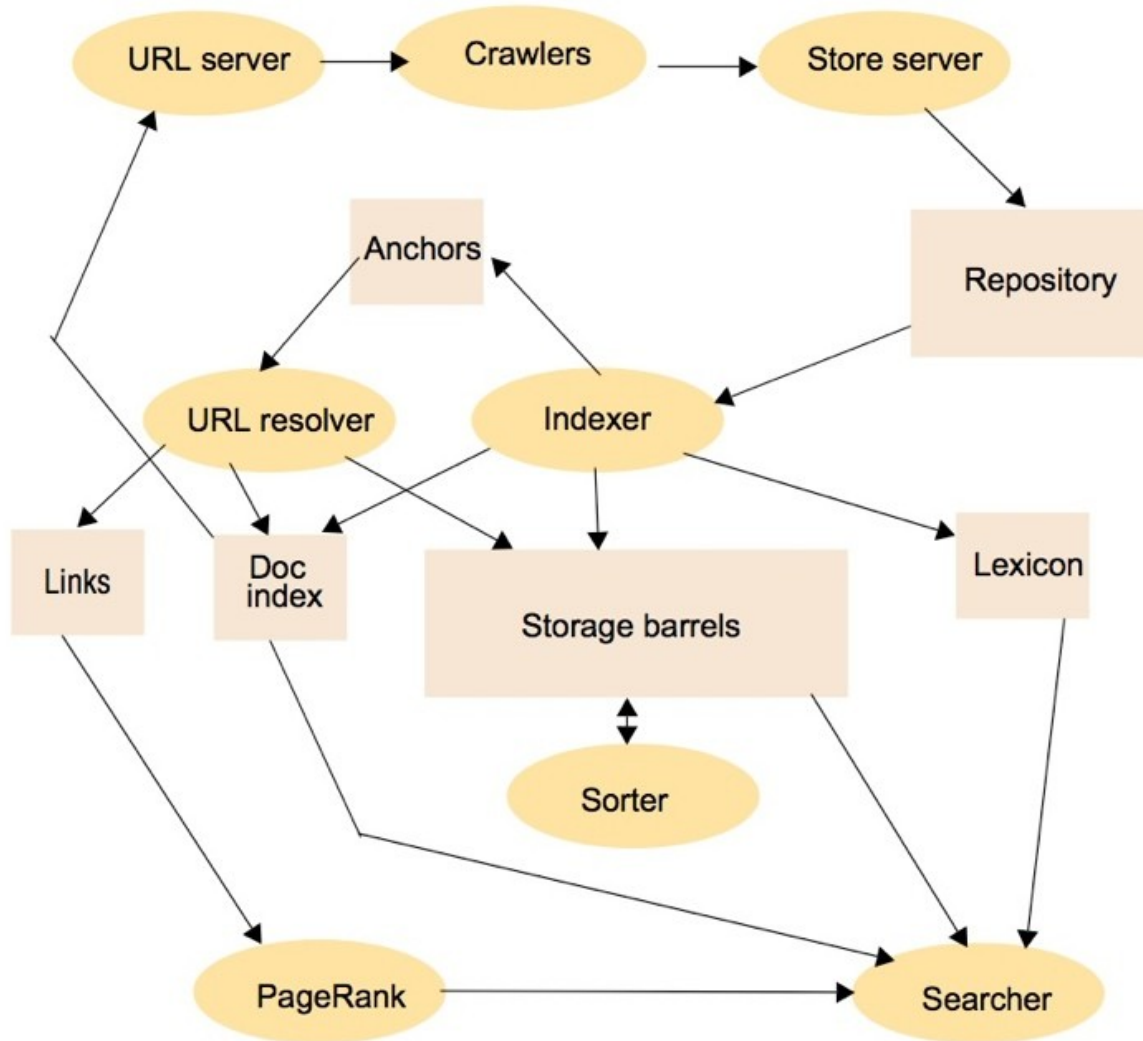


*From* **Coulouris, Dollimore, Kindberg and Blair**  
**Distributed Systems:  
Concepts and Design**

Edition 5, © Addison-Wesley 2012

Figure 21.1

Outline architecture of the original Google search engine [Brin and Page 1998]

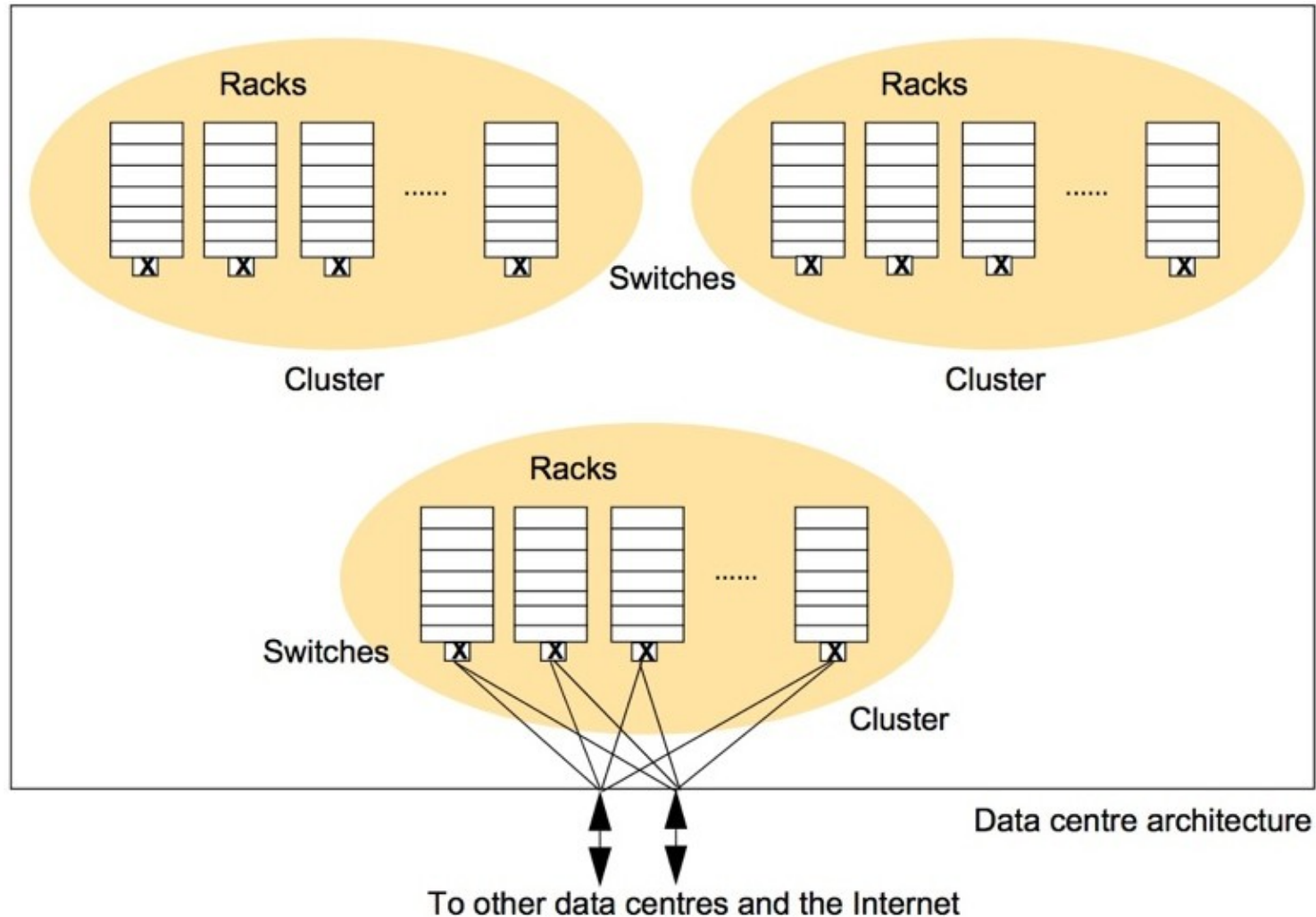


## Figure 21.2

### Example Google applications

| <i>Application</i> | <i>Description</i>  |
|--------------------|---|
| Gmail              | Mail system with messages hosted by Google but desktop-like message management.                           |
| Google Docs        | Web-based office suite supporting shared editing of documents held on Google servers.                     |
| Google Sites       | Wiki-like web sites with shared editing facilities.   |
| Google Talk        | Supports instant text messaging and Voice over IP.  |
| Google Calendar    | Web-based calendar with all data hosted on Google servers.  |
| Google Wave        | Collaboration tool integrating email, instant messaging, wikis and social networks.                       |
| Google News        | Fully automated news aggregator site.   |
| Google Maps        | Scalable web-based world map including high-resolution imagery and unlimited user-generated overlays.     |
| Google Earth       | Scalable near-3D view of the globe with unlimited user-generated overlays.                                |
| Google App Engine  | Google distributed infrastructure made available to outside parties as a service (platform as a service). |

Figure 21.3  
Organization of the Google physical infrastructure



(To avoid clutter the Ethernet connections are shown from only one of the clusters to the external links)

Figure 21.4  
The scalability problem in Google

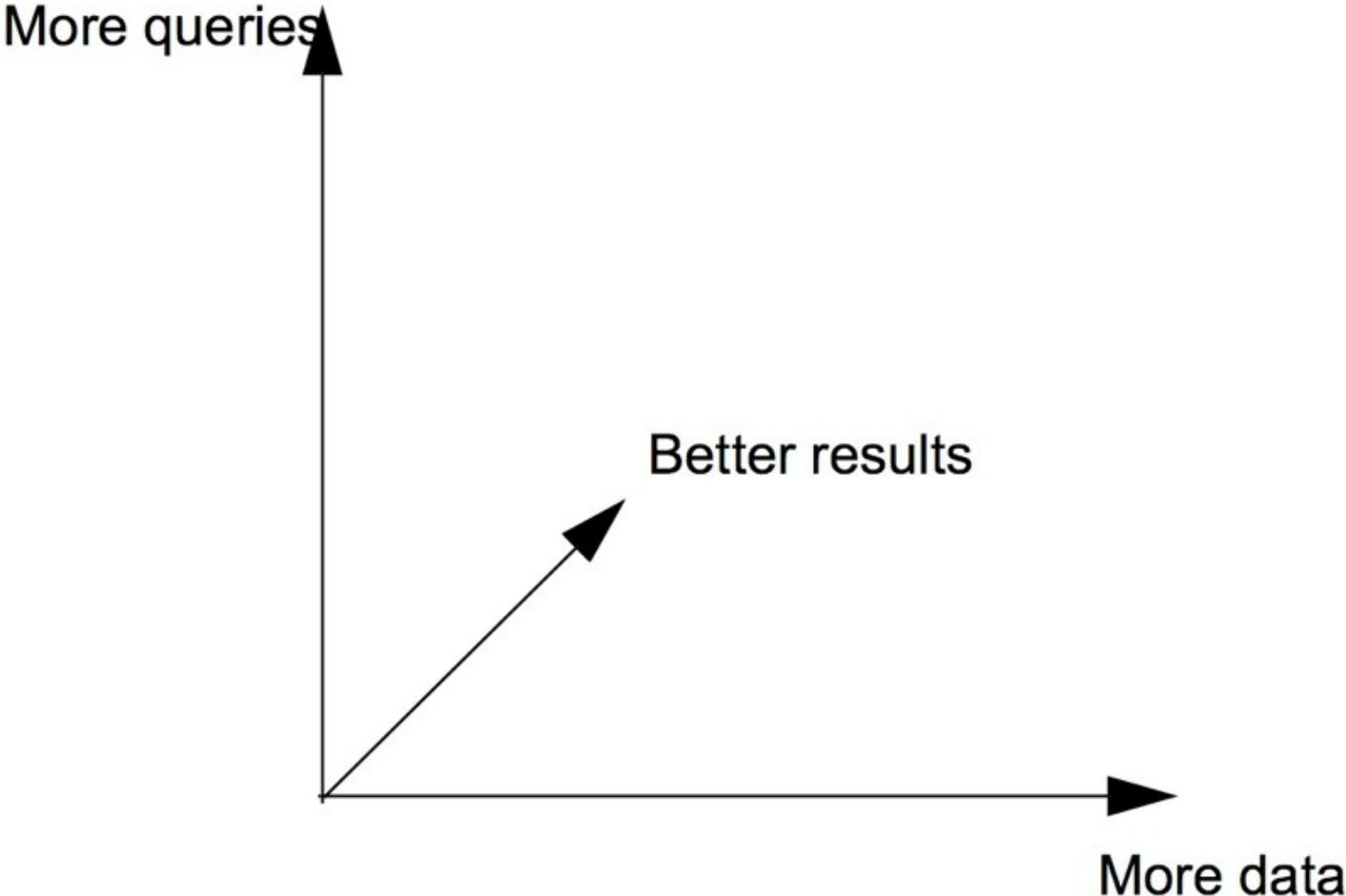


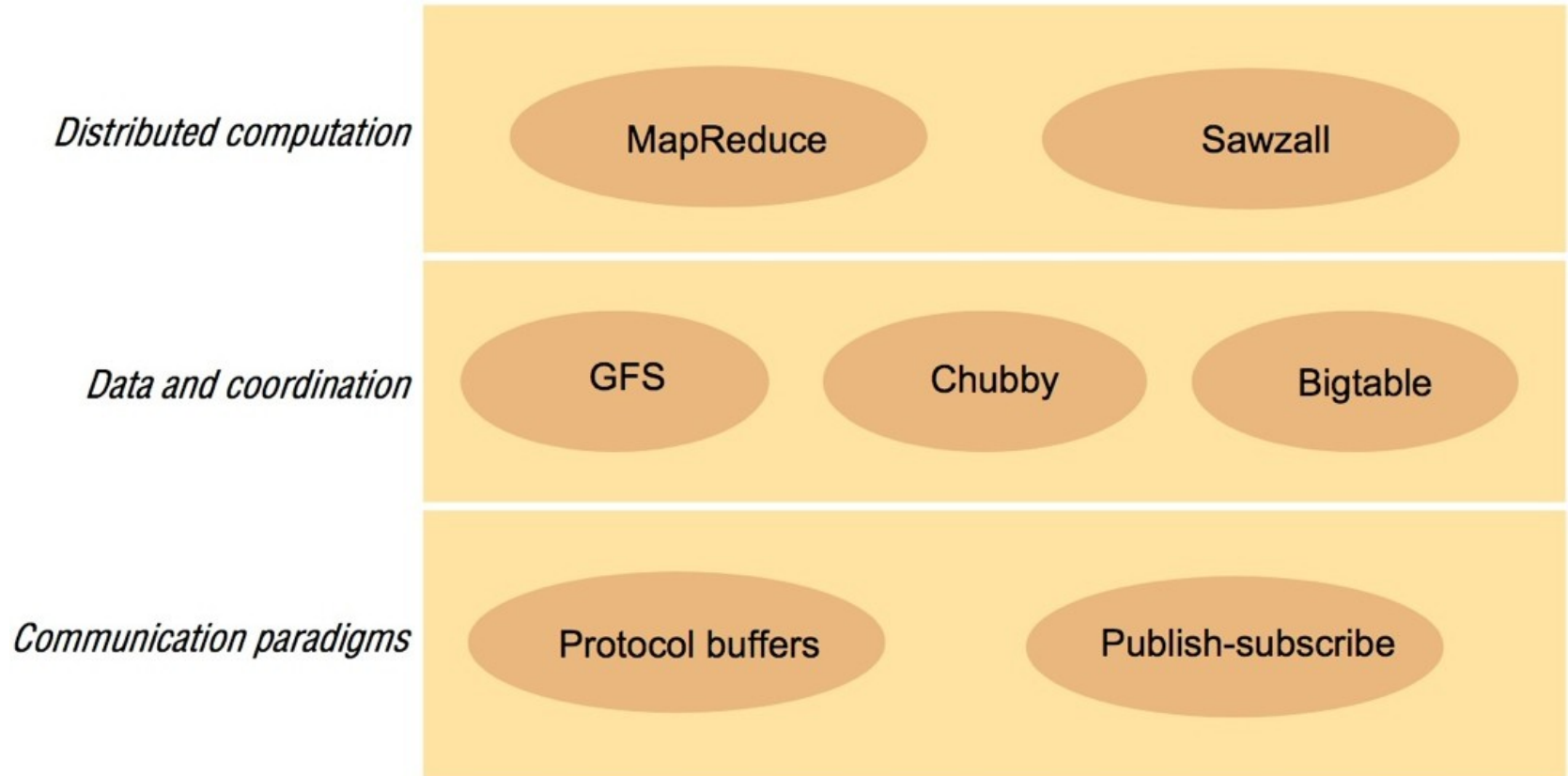
Figure 21.5  
The overall Google systems architecture

Google applications and services

Google infrastructure (middleware)

Google platform

Figure 21.6  
Google infrastructure



## Figure 21.7 Protocol buffers example

```
message Book {  
  required string title = 1;  
  repeated string author = 2;  
  enum Status {  
    IN_PRESS = 0;  
    PUBLISHED = 1;  
    OUT_OF_PRINT = 2;  
  }  
  message BookStats {  
    required int32 sales = 1;  
    optional int32 citations = 2;  
    optional Status bookstatus = 3 [default = PUBLISHED];  
  }  
  optional BookStats statistics = 3;  
  repeated string keyword = 4;  
}
```



## Figure 21.8a

### Summary of design choices related to communication paradigms - part 1

| <i>Element</i>   | <i>Design choice</i>   | <i>Rationale</i>   | <i>Trade-offs</i>   |
|------------------|--|--|---|
| Protocol buffers | The use of a language for specifying data formats  | Flexible in that the same language can be used for serializing data for storage or communication | -   |
|                  | Simplicity of the language   | Efficient implementation   | Lack of expressiveness when compared, for example, with XML         |
|                  | Support for a style of RPC (taking a single message as a parameter and returning a single message as result) | More efficient, extensible and supports service evolution  | Lack of expressiveness when compared with other RPC or RMI packages |
|                  | Protocol-agnostic design   | Different RPC implementations can be used  | No common semantics for RPC exchanges                               |

## Figure 21.8b

### Summary of design choices related to communication paradigms - *part 2*

---

|                   |                                      |   |  |
|-------------------|--------------------------------------|---|--|
| Publish-subscribe | Topic-based approach                 | Supports efficient implementation                           | Less expressive than content-based approaches (mitigated by the additional filtering capabilities) |
|                   | Real-time and reliability guarantees | Supports maintenance of consistent views in a timely manner | Additional algorithmic support required with associated overhead                                   |

---

Figure 21.9  
Overall architecture of GFS

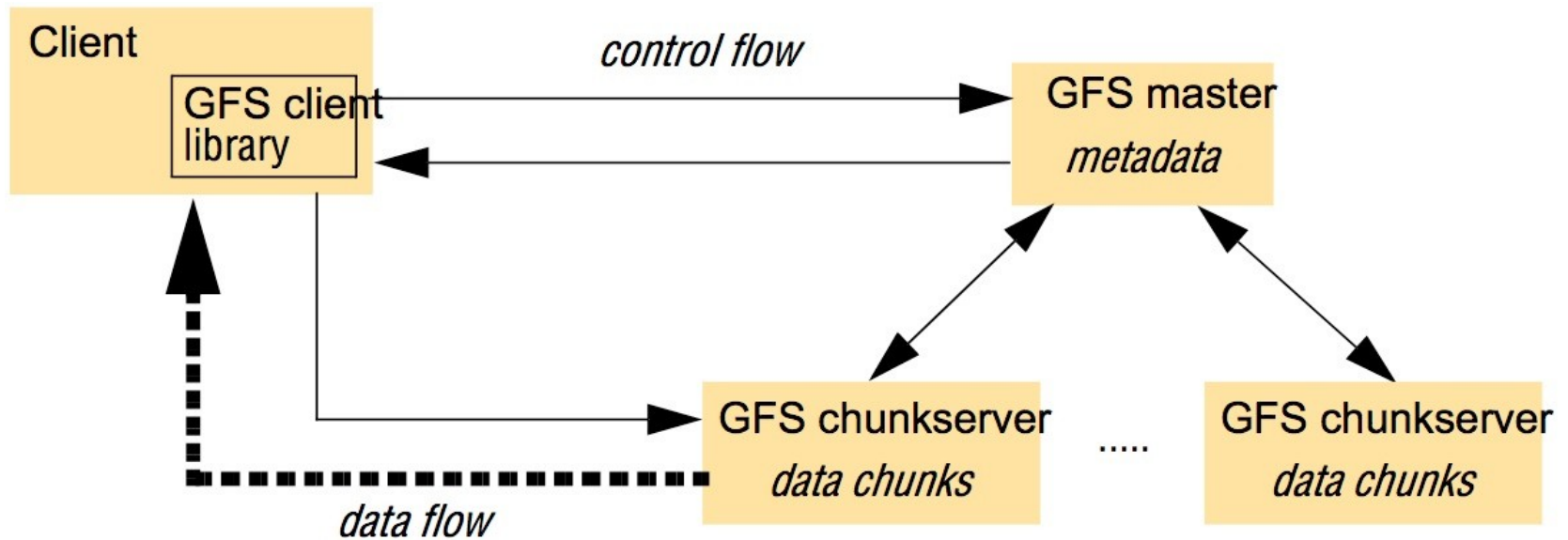


Figure 21.10  
Chubby API

| <i>Role</i> | <i>Operation</i>          | <i>Effect</i>   |
|-------------|---------------------------|---|
| General     | <i>Open</i>               | Opens a given named file or directory and returns a handle                            |
|             | <i>Close</i>              | Closes the file associated with the handle  |
|             | <i>Delete</i>             | Deletes the file or directory   |
| File        | <i>GetContentsAndStat</i> | Returns (atomically) the whole file contents and metadata associated with the file    |
|             | <i>GetStat</i>            | Returns just the metadata   |
|             | <i>ReadDir</i>            | Returns the contents of a directory – that is, the names and metadata of any children |
|             | <i>SetContents</i>        | Writes the whole contents of a file (atomically)                                      |
|             | <i>SetACL</i>             | Writes new access control list information  |
| Lock        | <i>Acquire</i>            | Acquires a lock on a file   |
|             | <i>TryAcquire</i>         | Tries to acquire a lock on a file   |
|             | <i>Release</i>            | Releases a lock   |

Figure 21.11  
Overall architecture of Chubby

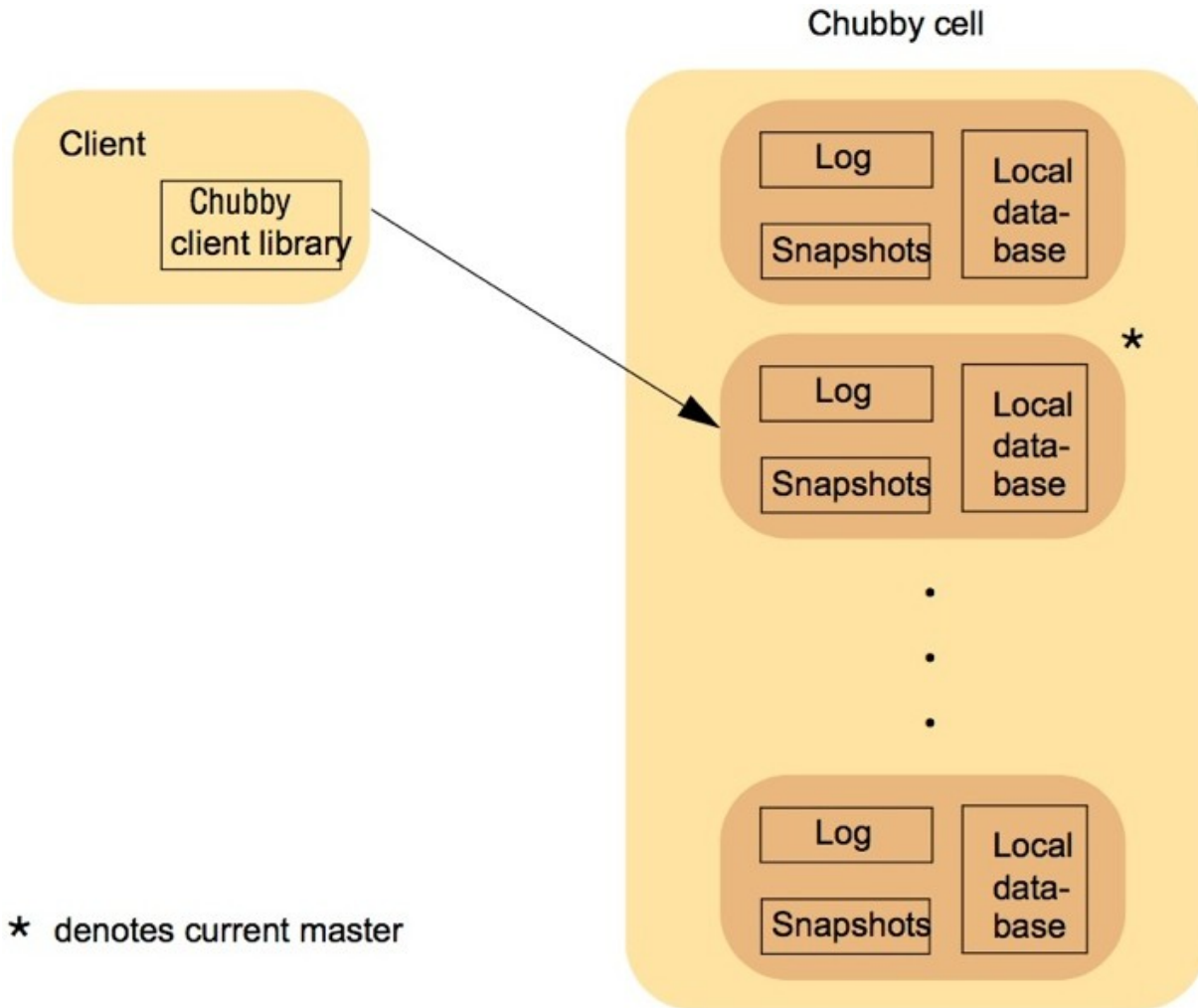


Figure 21.12  
Message exchanges in Paxos (in absence of failures) - step 1

Step 1: electing a coordinator

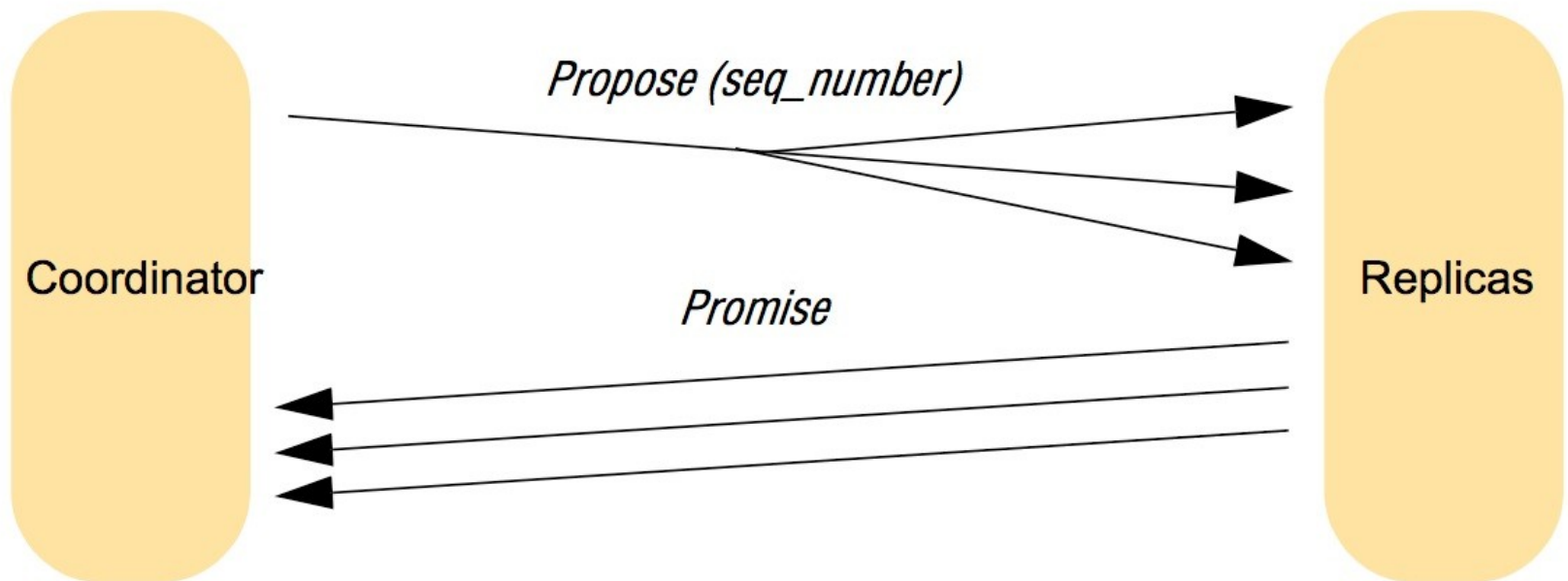


Figure 21.12  
Message exchanges in Paxos (in absence of failures) - step 2

Step 2: seeking consensus

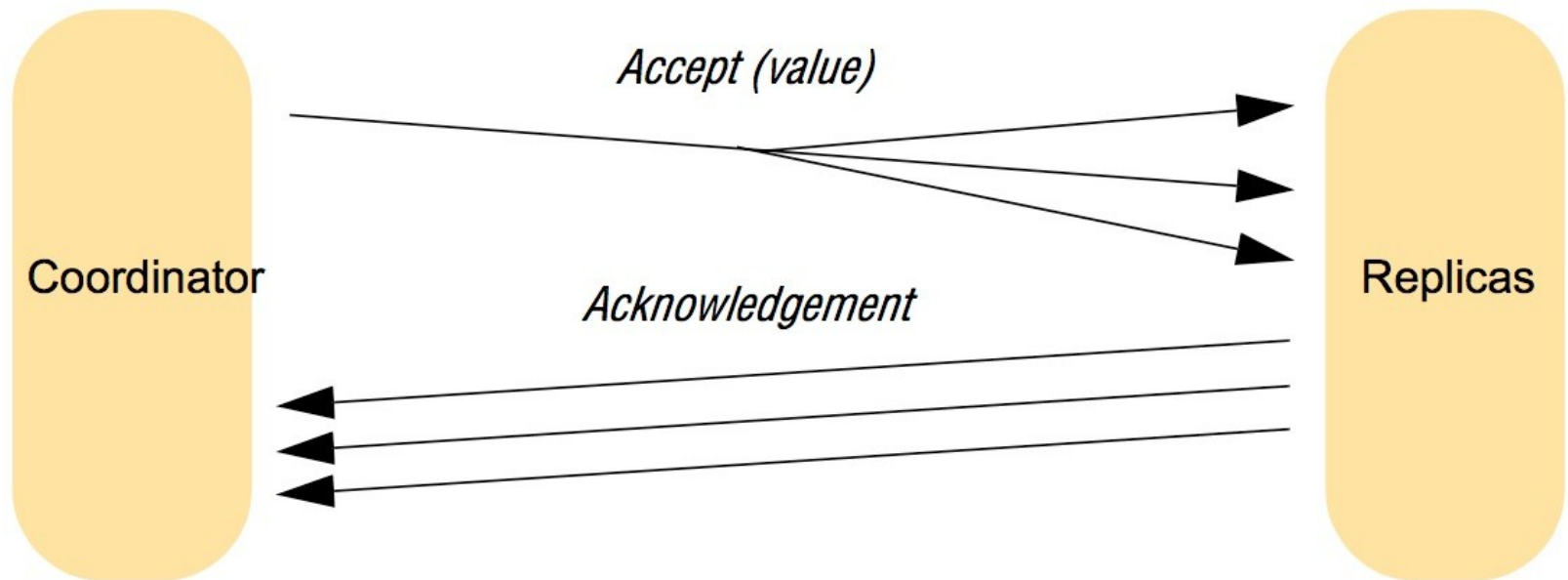


Figure 21.12  
Message exchanges in Paxos (in absence of failures) - step 3

---

Step 3: achieving consensus

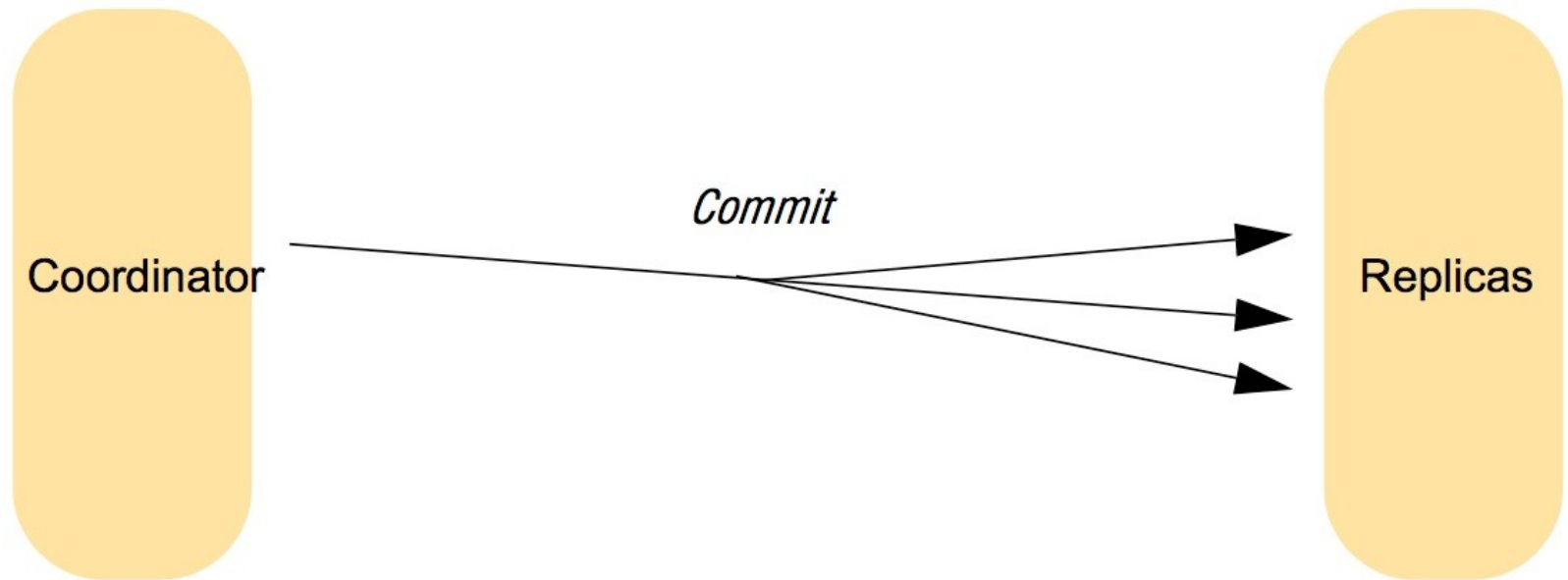




Figure 21.13  
The table abstraction in Bigtable

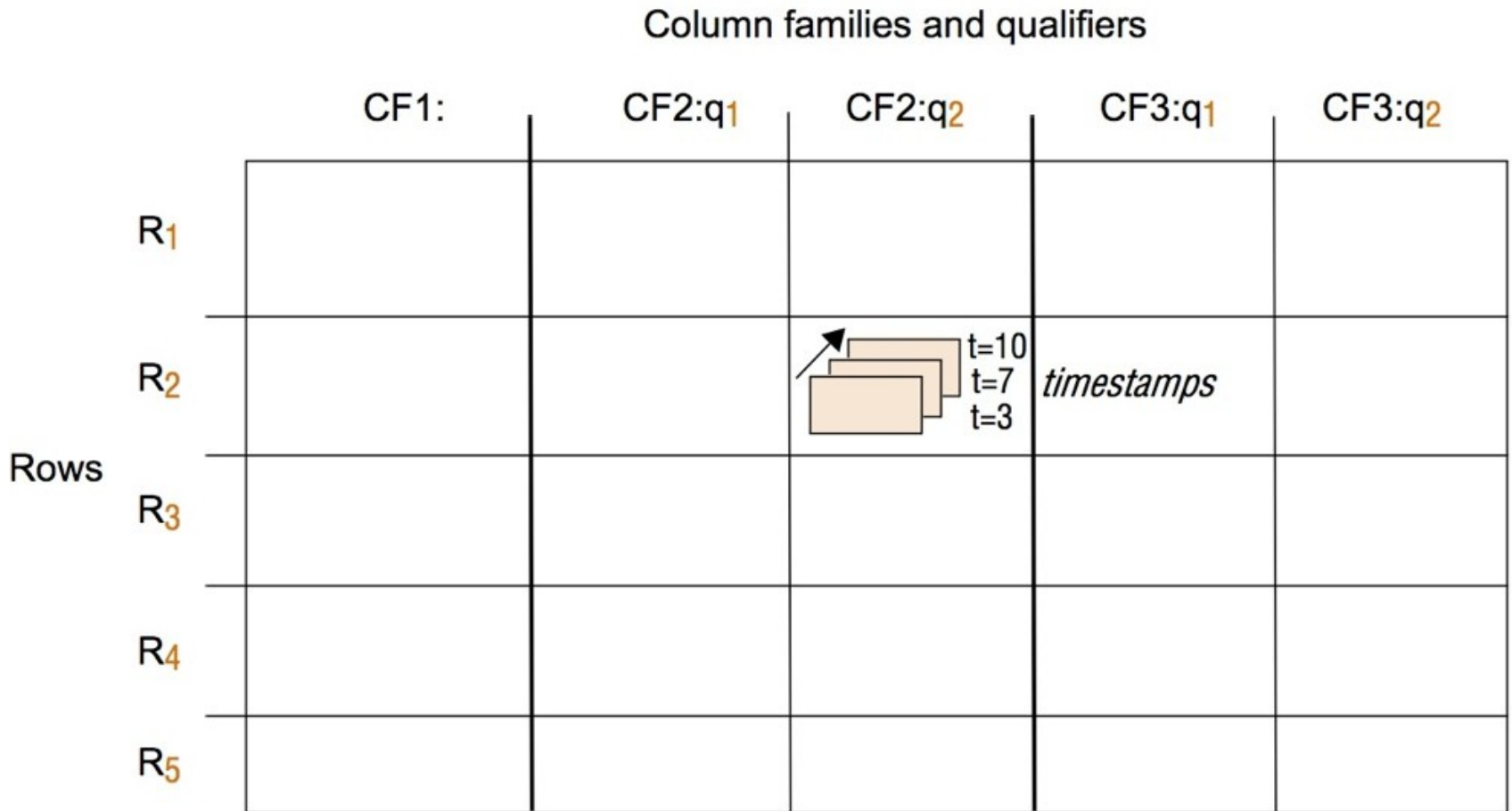


Figure 21.14  
Overall architecture of Bigtable

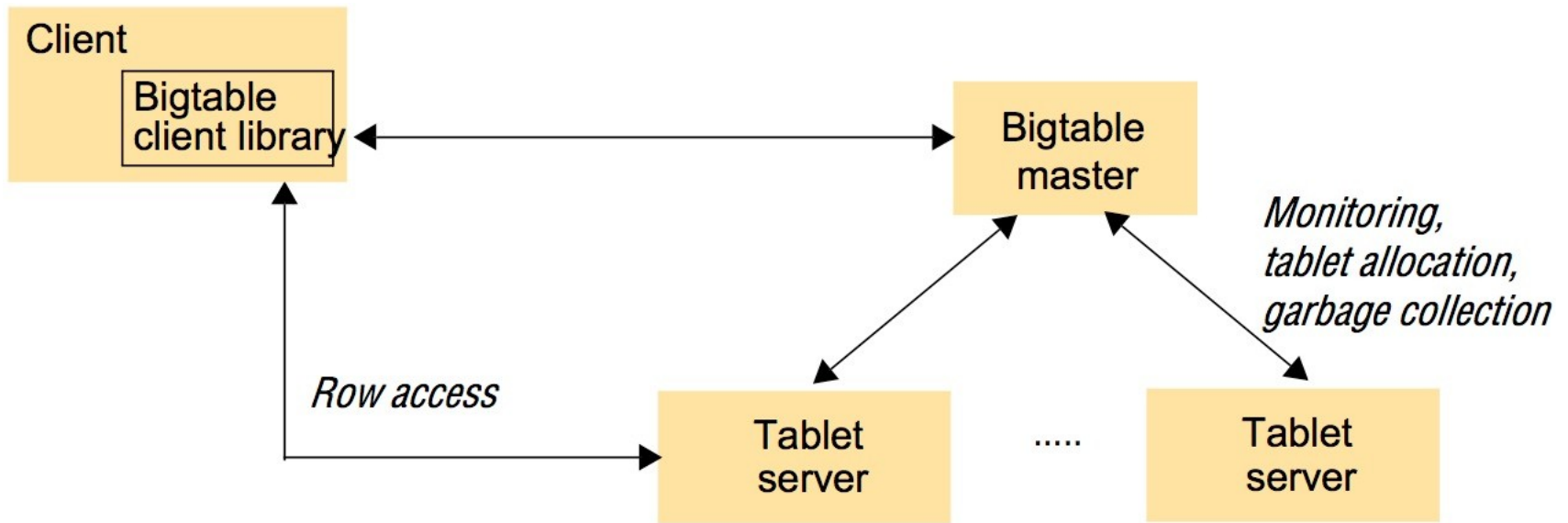


Figure 21.15  
The storage architecture in Bigtable

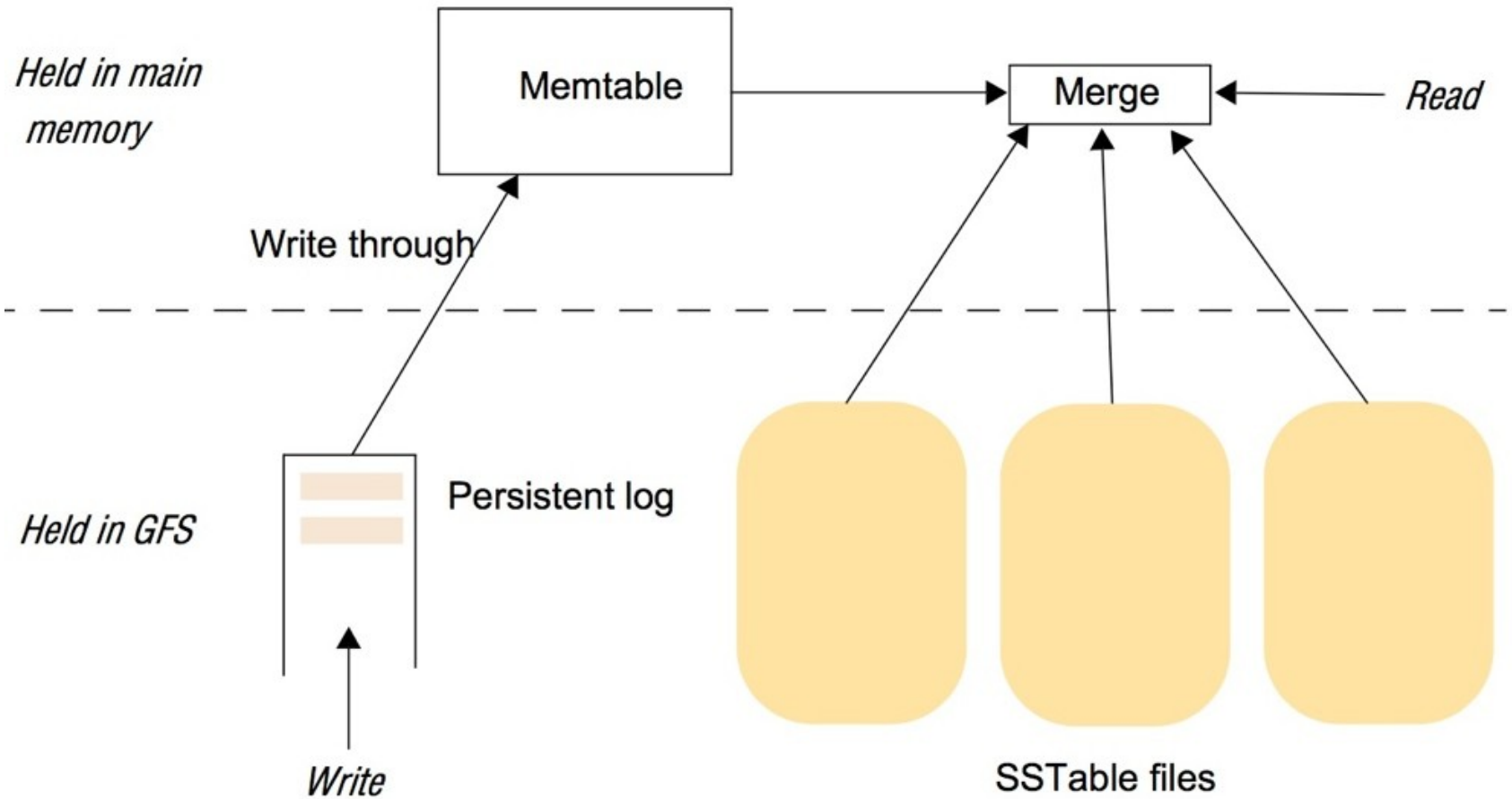


Figure 21.16  
The hierarchical indexing scheme adopted by Bigtable

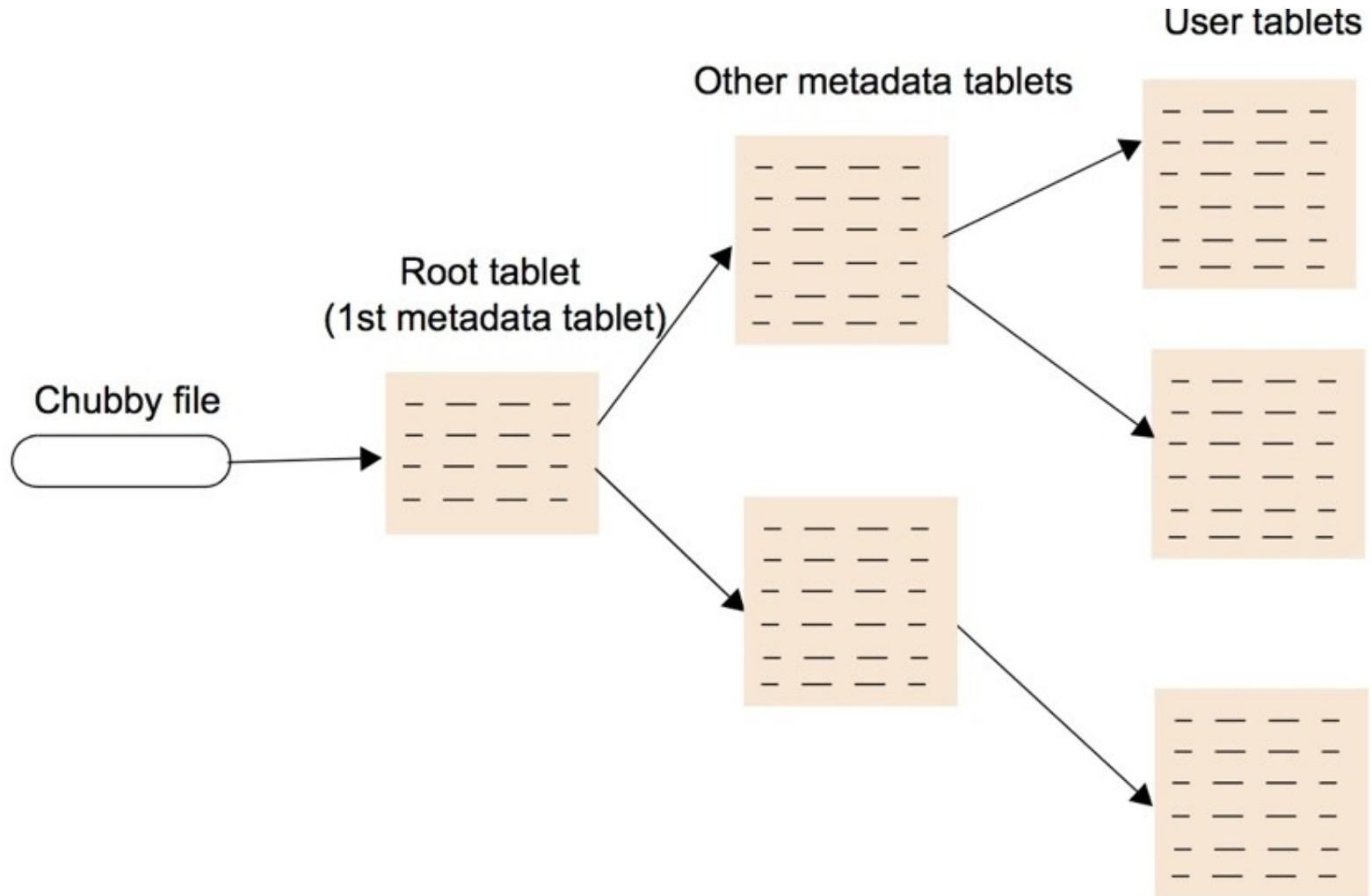


Figure 21.17  
Summary of design choices related to data storage and coordination

| <i>Element</i>  | <i>Design choice</i>                         | <i>Rationale</i>   | <i>Trade-offs</i>  |
|-----------------|--|--|--|
| <i>GFS</i>      | The use of a large chunk size (64 megabytes) | Suited to the size of files in GFS; efficient for large sequential reads and appends; minimizes the amount of metadata | Would be very inefficient for random access to small parts of files                  |
|                 | The use of a centralized master              | The master maintains a global view that informs management decisions; simpler to implement                             | Single point of failure (mitigated by maintaining replicas of operations logs)       |
|                 | Separation of control and data flows         | High-performance file access with minimal master involvement   | Complicates the client library as it must deal with both the master and chunkservers |
|                 | Relaxed consistency model                    | High performance, exploiting semantics of the GFS operations   | Data may be inconsistent, in particular duplicated                                   |
| <i>Chubby</i>   | Combined lock and file abstraction           | Multipurpose, for example supporting elections   | Need to understand and differentiate between different facets                        |
|                 | Whole-file reading and writing               | Very efficient for small files   | Inappropriate for large files  |
|                 | Client caching with strict consistency       | Deterministic semantics  | Overhead of maintaining strict consistency   |
| <i>Bigtable</i> | The use of a table abstraction               | Supports structured data efficiently   | Less expressive than a relational database   |
|                 | The use of a centralized master              | As above, master has a global view; simpler to implement   | Single point of failure; possible bottleneck   |
|                 | Separation of control and data flows         | High-performance data access with minimal master involvement   | -  |
|                 | Emphasis on monitoring and load balancing    | Ability to support very large numbers of parallel clients  | Overhead associated with maintaining global states                                   |

Figure 21.18  
Examples of the use of MapReduce

| <i>Function</i>   | <i>Initial step</i>   | <i>Map phase</i>   | <i>Intermediate step</i>   | <i>Reduce phase</i>   |
|---|---|--|--|---|
| <i>Word count</i>   | <i>Partition data into fixed-size chunks for processing</i> | For each occurrence of word in data partition, emit $\langle \text{word}, 1 \rangle$   | <i>Merge/sort all key-value keys according to their intermediary key</i> | For each word in the intermediary set, count the number of 1s |
| <i>Grep</i>   |   | Output a line if it matches a given pattern  |  | Null  |
| <i>Sort</i><br><i>N.B. This relies heavily on the intermediate step</i> |   | For each entry in the input data, output the key-value pairs to be sorted  |  | Null  |
| <i>Inverted index</i>   |   | Parse the associated documents and output a $\langle \text{word}, \text{document ID} \rangle$ pair wherever that word exists |  | For each word, produce a list of (sorted) document IDs        |

Figure 21.19  
The overall execution of a MapReduce program

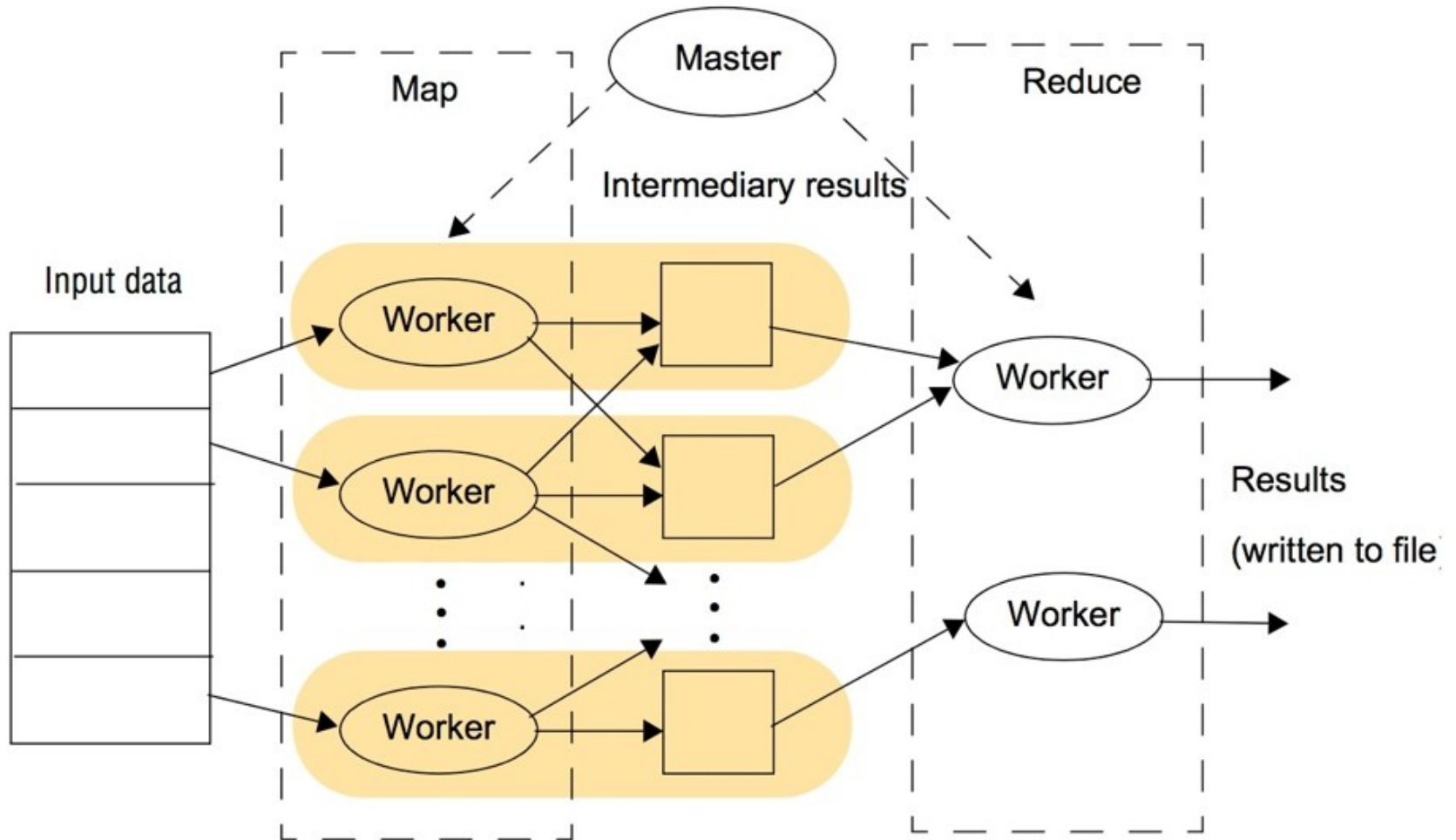
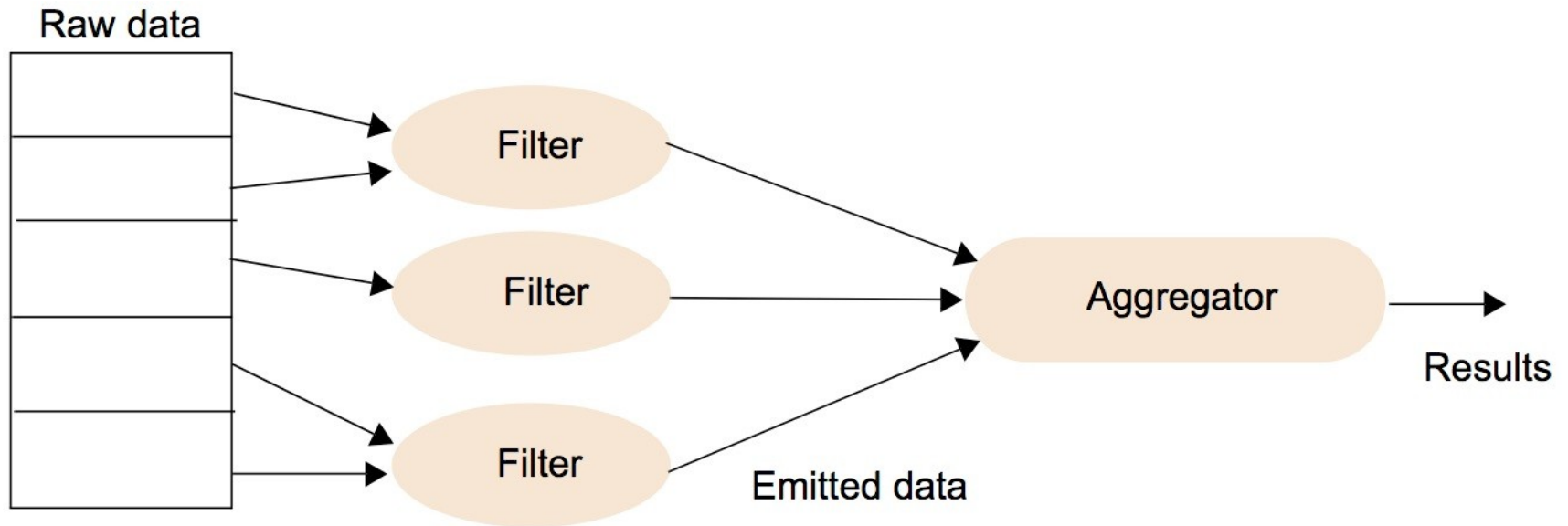




Figure 21.20  
The overall execution of a Sawzall program





## Figure 21.21 Summary of design choices related to distributed computation

| <i>Element</i>   | <i>Design choice</i>  | <i>Rationale</i>   | <i>Trade-offs</i>  |
|------------------|---|--|--|
| <i>MapReduce</i> | The use of a common framework   | Hides details of parallelization and distribution from the programmer; improvements to the infrastructure immediately exploited by all MapReduce applications                  | Design choices within the framework may not be appropriate for all styles of distributed computation |
|                  | Programming of system via two operations, <i>map</i> and <i>reduce</i>      | Very simple programming model allowing rapid development of complex distributed computations   | Again, may not be appropriate for all problem domains  |
|                  | Inherent support for fault-tolerant distributed computations                | Programmer does not need to worry about dealing with faults (particularly important for long-running tasks running over a physical infrastructure where failures are expected) | Overhead associated with fault-recovery strategies   |
| <i>Sawzall</i>   | Provision of a specialized programming language for distributed computation | Again, support for rapid development of often complex distributed computations with complexity hidden from the programmer (even more so than with MapReduce)                   | Assumes that programs can be written in the style supported (in terms of filters and aggregators)    |