

# GIGA+: Scalable Directories for Shared File Systems (or, How to build directories with trillions of files)

Swapnil V. Patil and Garth A. Gibson

## Problem: Scalable Directories

Need high performance metadata services

- Most file systems store a directory on a single MDS
- New trends need large metadata services
  - Apps generating millions of small files in a directory, like a simple database
  - Large apps run in parallel on clusters of 100,000s of CPUs

Build scalable directories for shared file-systems

- POSIX-compliant, maintain UNIX file system semantics
- Store trillions of files and handle >100K operations/second

## Goal: More Scalability Through More Parallelism

Minimize serialization

- Avoid ordered splitting of partitions, like LH\* [Litwin96]

Eliminate system-wide synchronization

- Avoid using cache consistency and distributed locking, like GPFS [Schmuck02]

GIGA+ distributed indexing divides a directory into partitions, spread across multiple servers

- Enables highly incremental, unsynchronized, and load-balanced growth

## GIGA+ Technique

Allows servers to grow their partitions independently

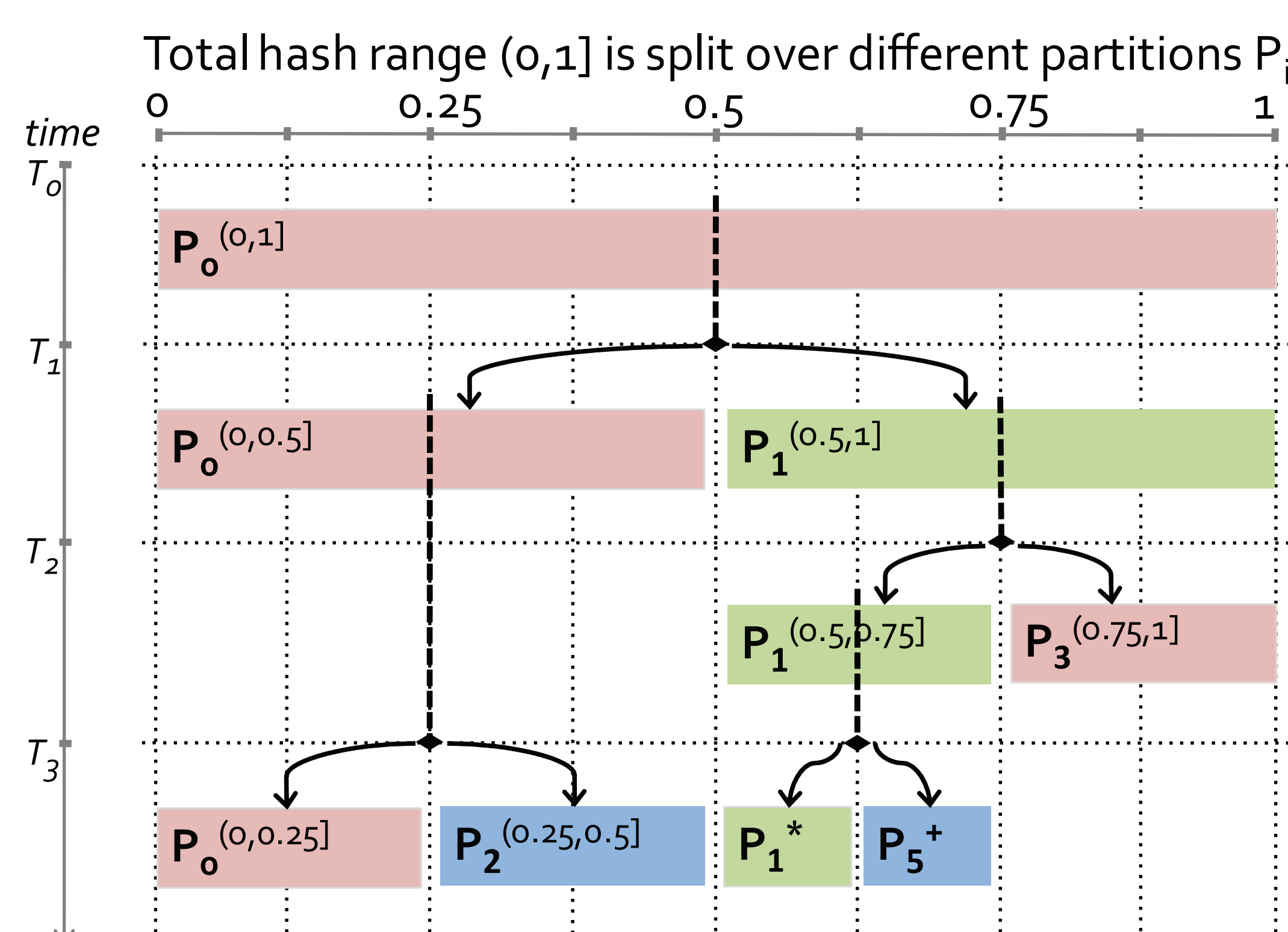
- Only maintain local state about their partitions
- Keep "split history" of their partitions

Tolerates out-of-date partition-to-server maps at the client

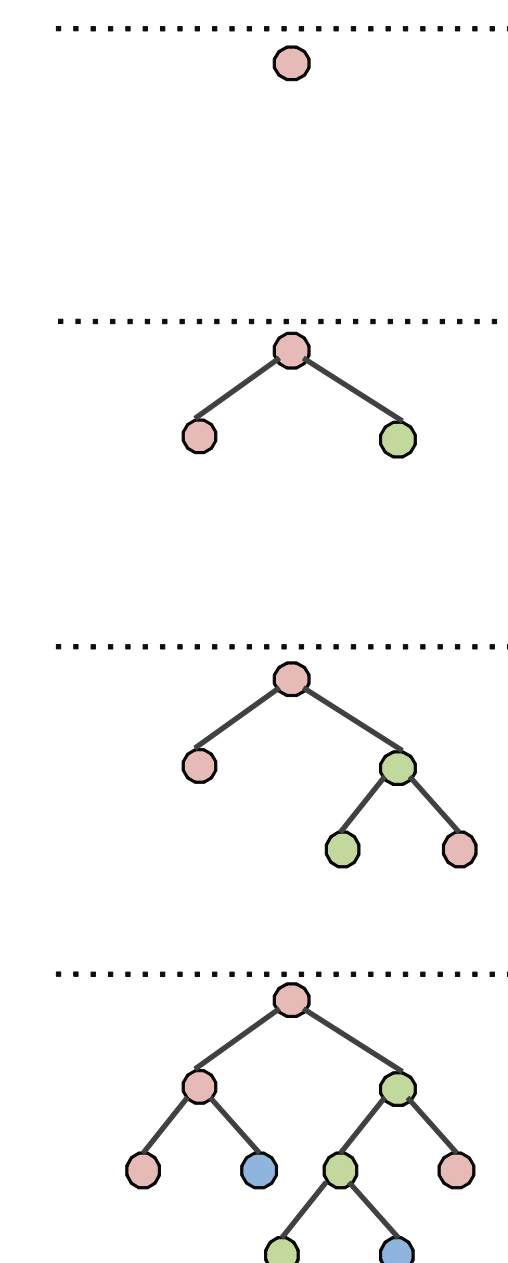
- Due to unsynchronized growth, map becomes state & inconsistent
- Copies updated lazily, on addressing an incorrect server

Unique, self-describing bitmap to map partitions on a server

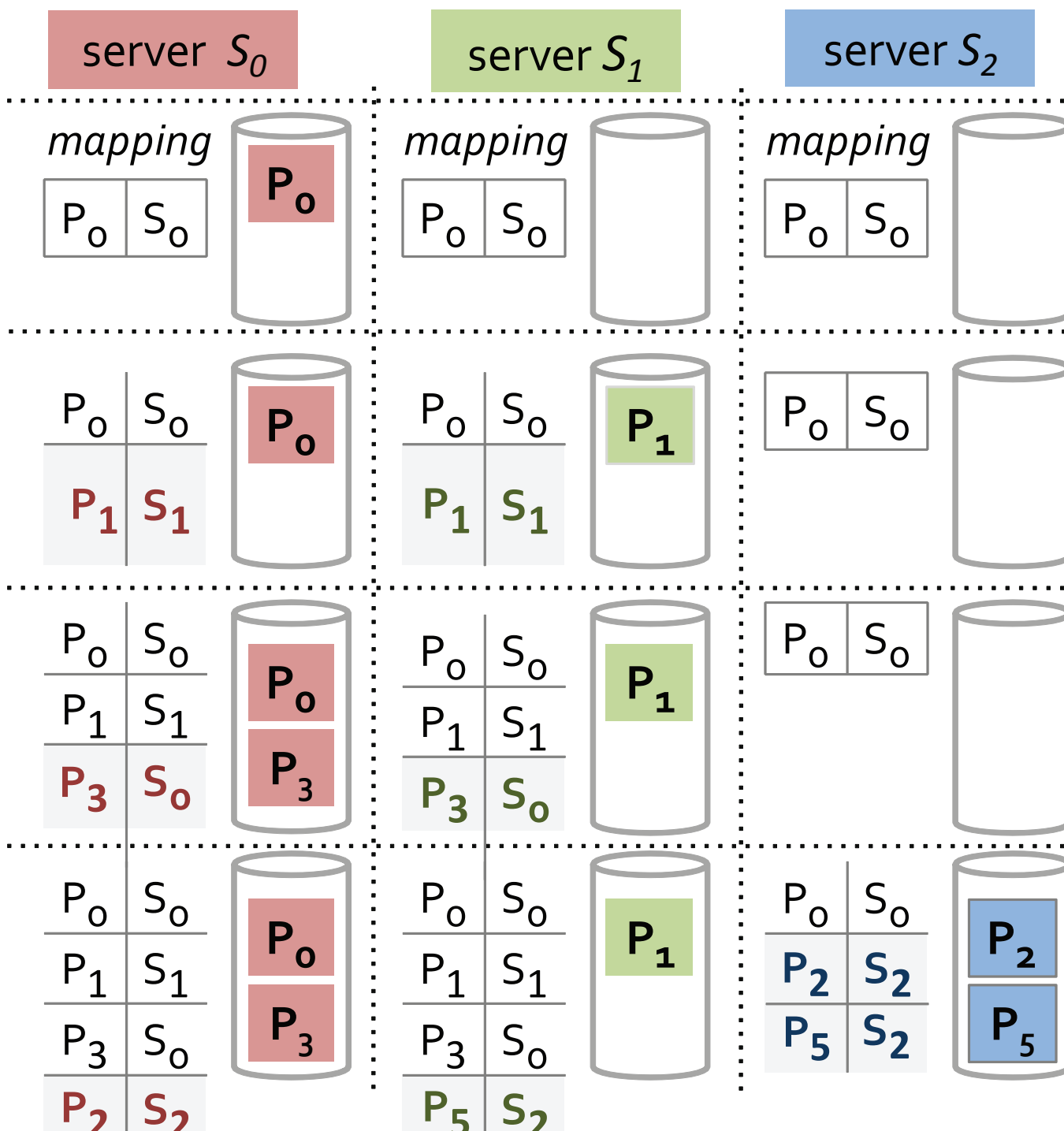
- Tracks presence or absence of a partition and its split history
- Deterministic search of best server to send the request
- Efficiently send many bitmap updates to erroneous clients
- Compact: billion file directory, in 16 KB



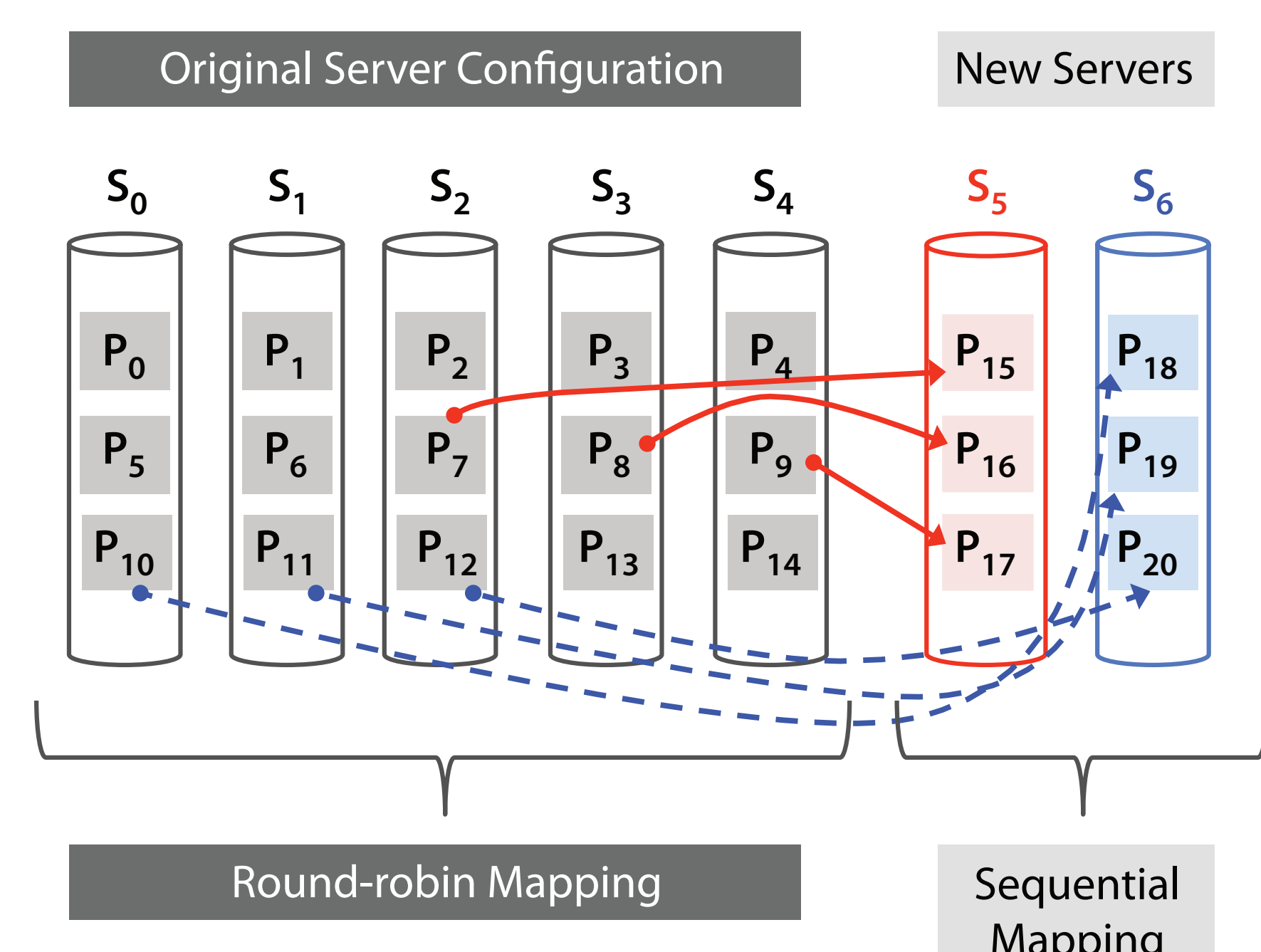
Logical view of the entire index



Physical view (mapping & partitions) on each server



## Reconfiguration and Recovering in GIGA+



Handling server addition

- Change the partition-to-server mapping from round-robin on the original server set to sequential on the newly added servers
- Minimizes the amount of data migration during reconfiguration

Handling failures

- Servers use "uniform de-clustered replication" that deterministically replicates each server's state spread across all remaining servers
- Enables load-balanced failover and fast, parallel recovery

