**XTREEMFS**

# BabuDB: Fast and Efficient File System Metadata Storage

*Jan Stender, Björn Kolbeck, Mikael Högqvist*

*Felix Hupfeld*

Zuse Institute Berlin

Google GmbH Zurich

## Motivation

- Modern parallel / distributed file systems:

  - Huge numbers of files and directories

  - Many storage servers but few metadata servers

- Examples:

  - Lustre, Panasas Active Scale, Google File System

- Metadata access critical wrt. system performance

  - ~75% of all file system calls are metadata accesses
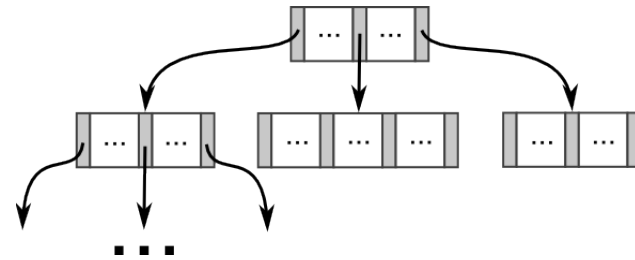
  - Metadata servers are bottlenecks

- B-tree-like data structures used for metadata storage
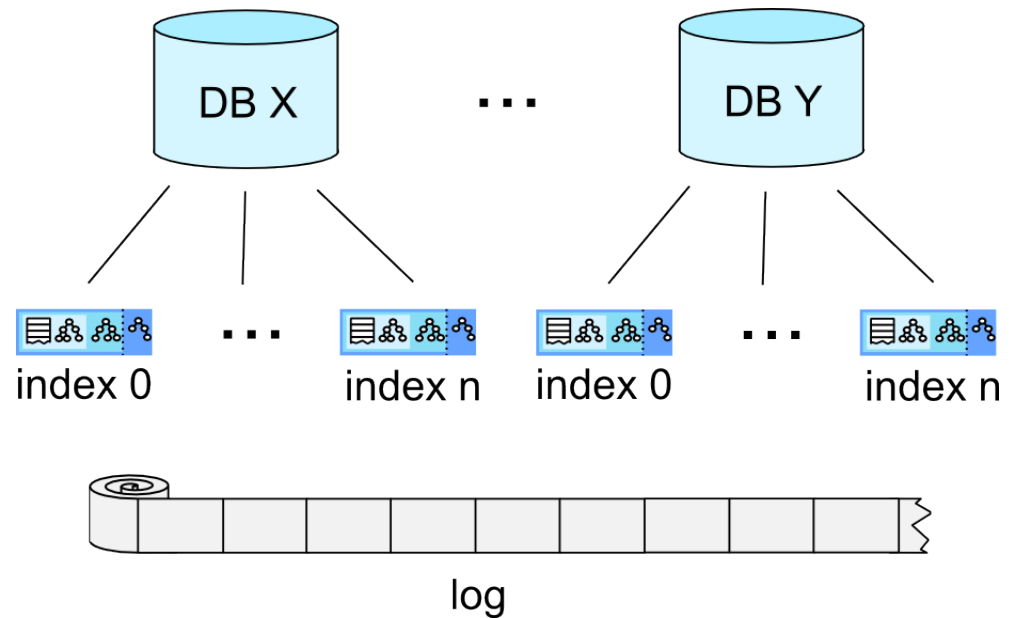
    - ZFS, btrfs, Lustre, PVFS2

- Downsides:

    - Hard to implement and test, high code complexity

    - Multi-version B-trees even more complex
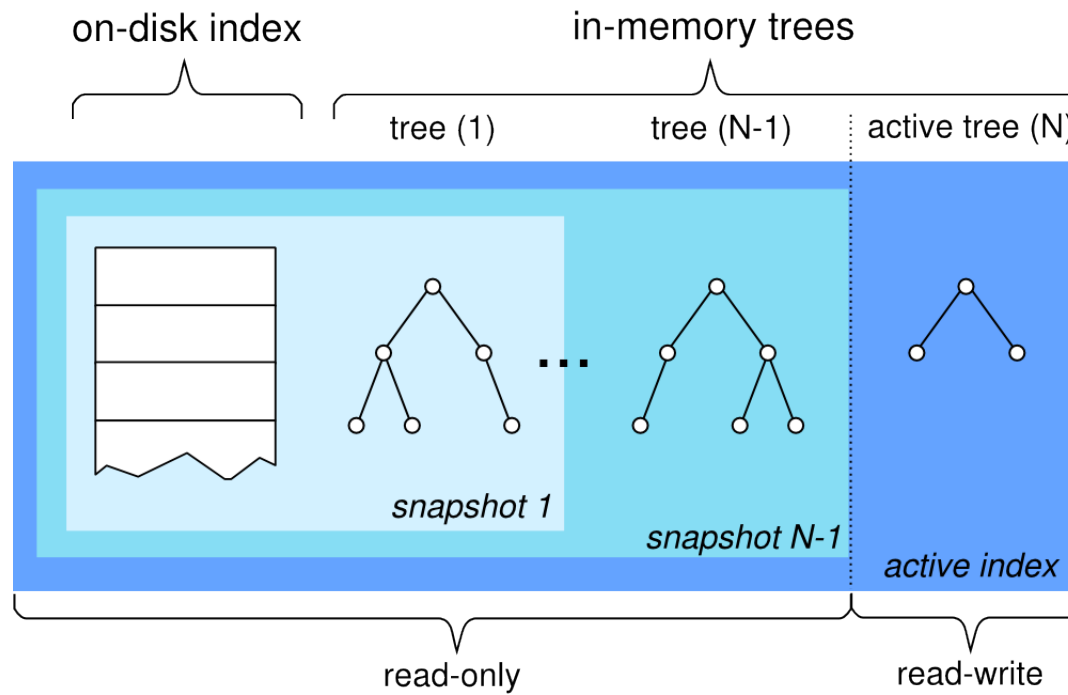
    - On-disk re-balancing expensive
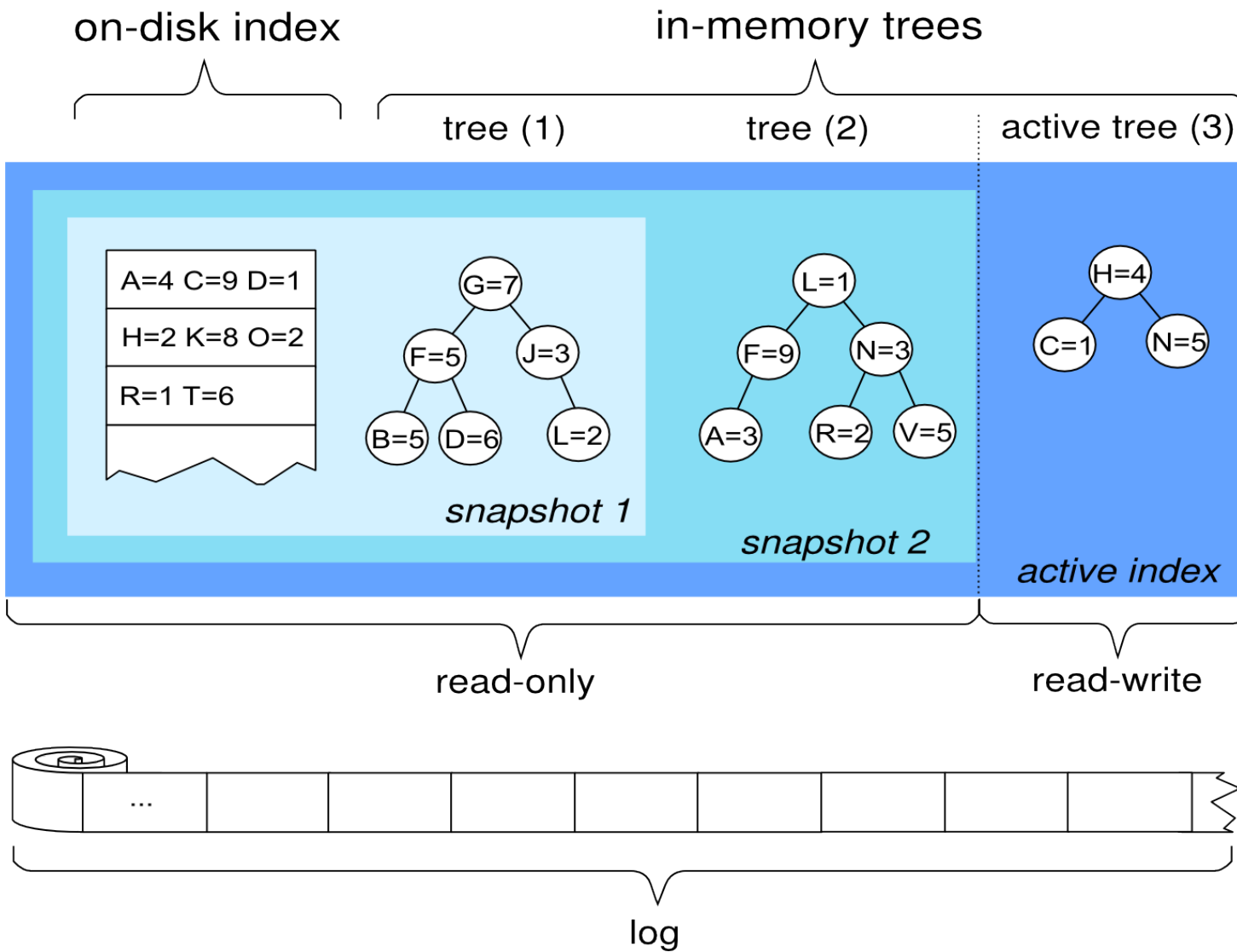
# BabuDB

– ## Key-value store



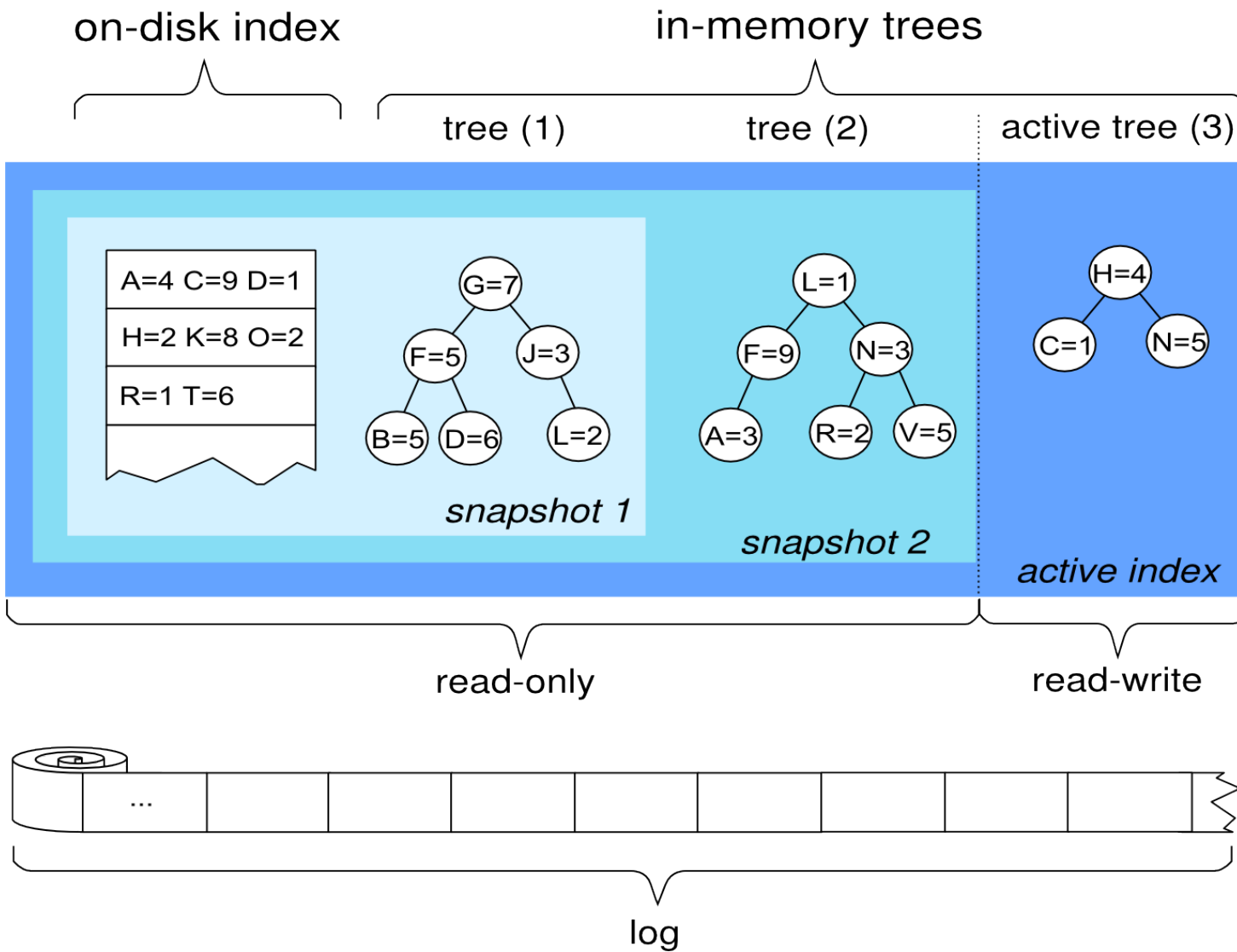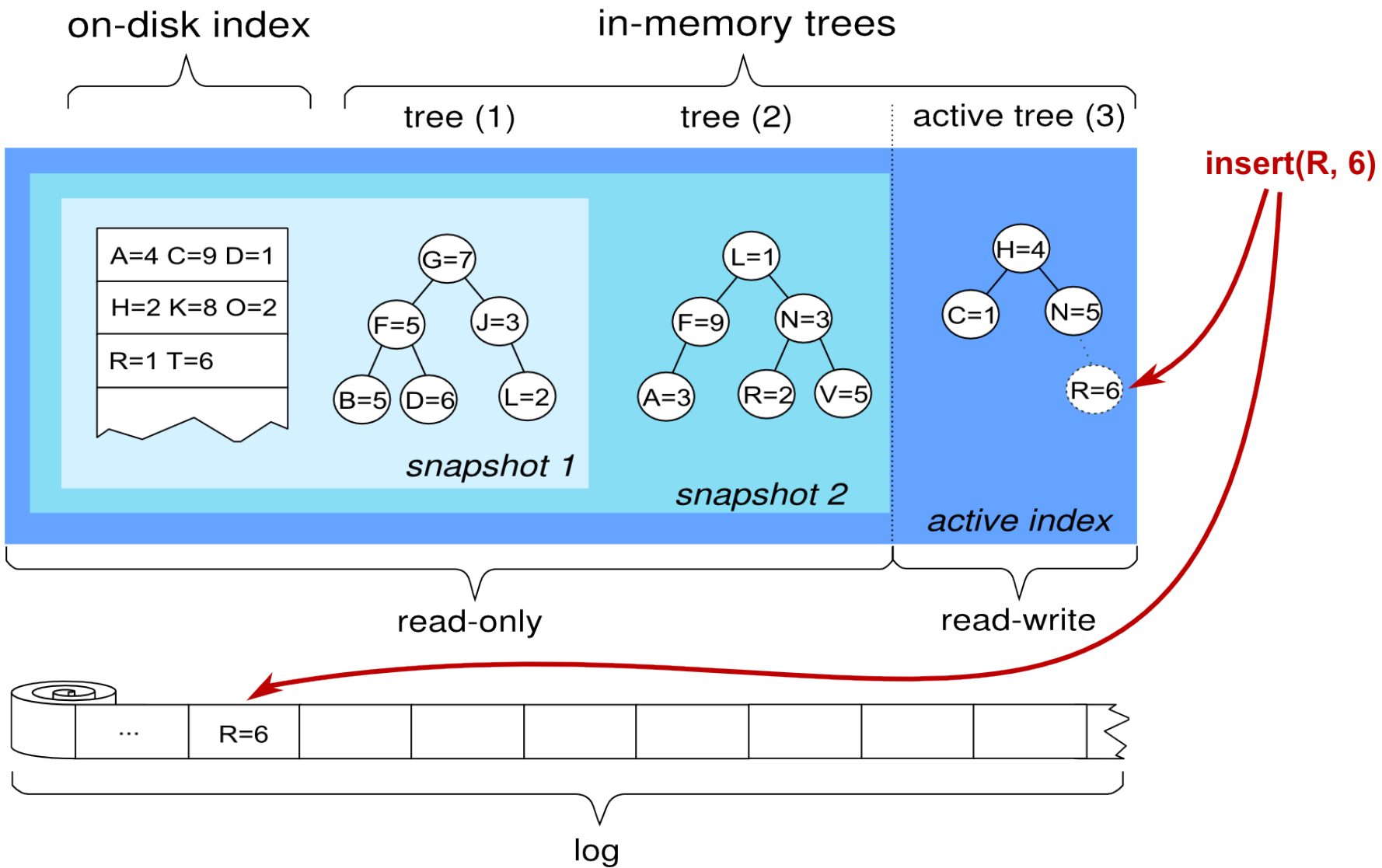– ## FS metadata: key-value pairs stored in DB indices

# BabuDB: Index

# Example

# Example: Insertions
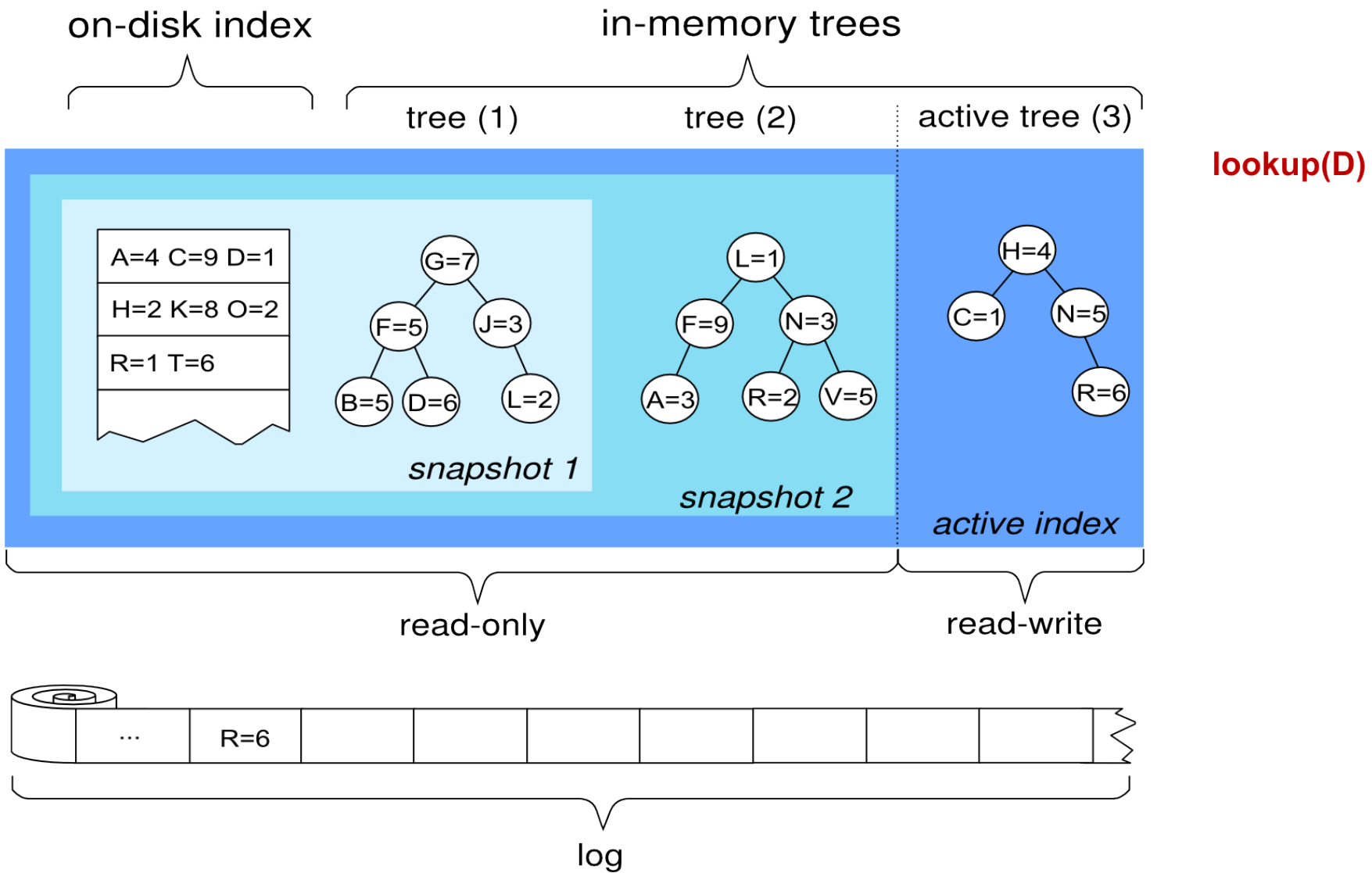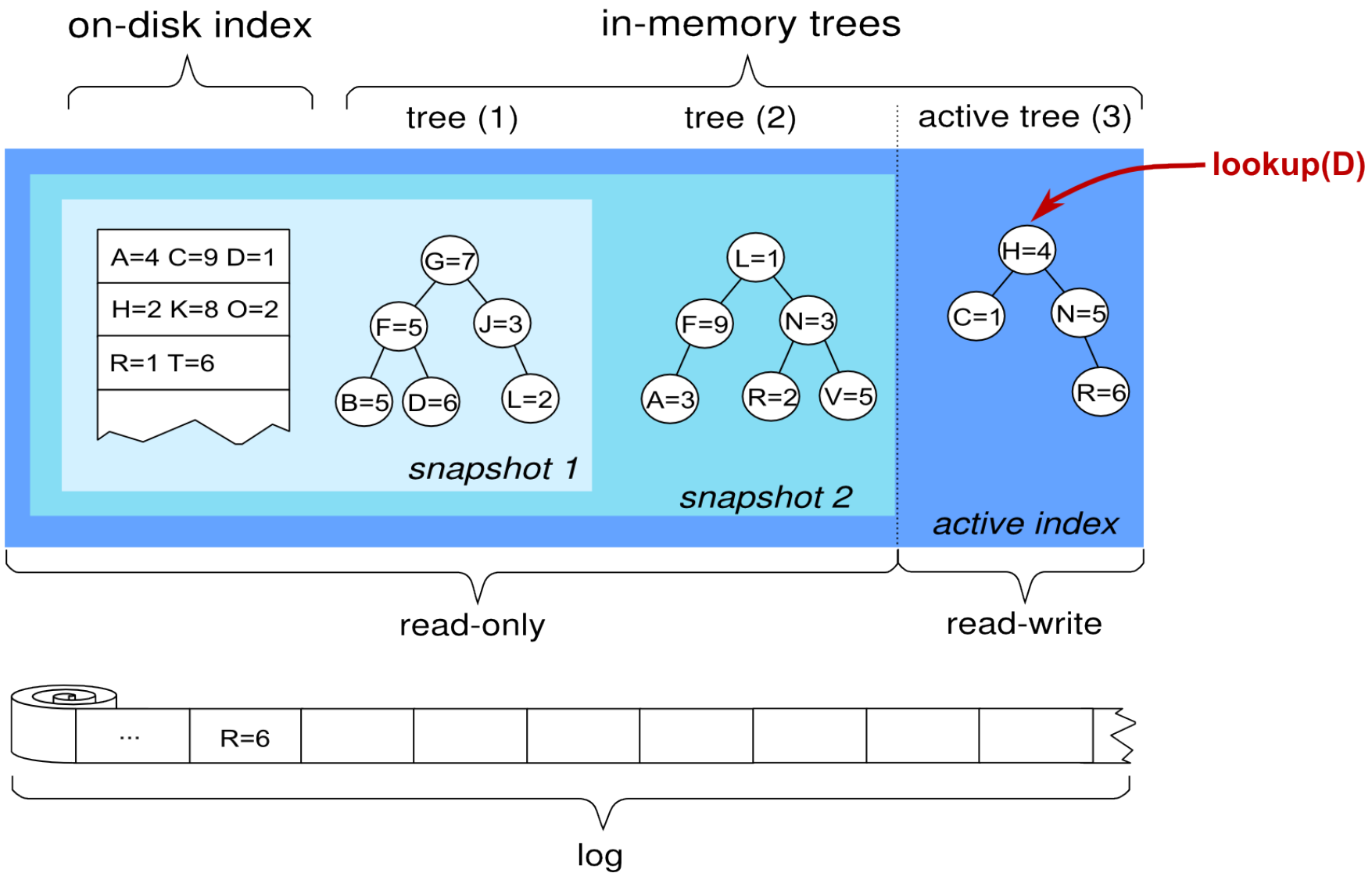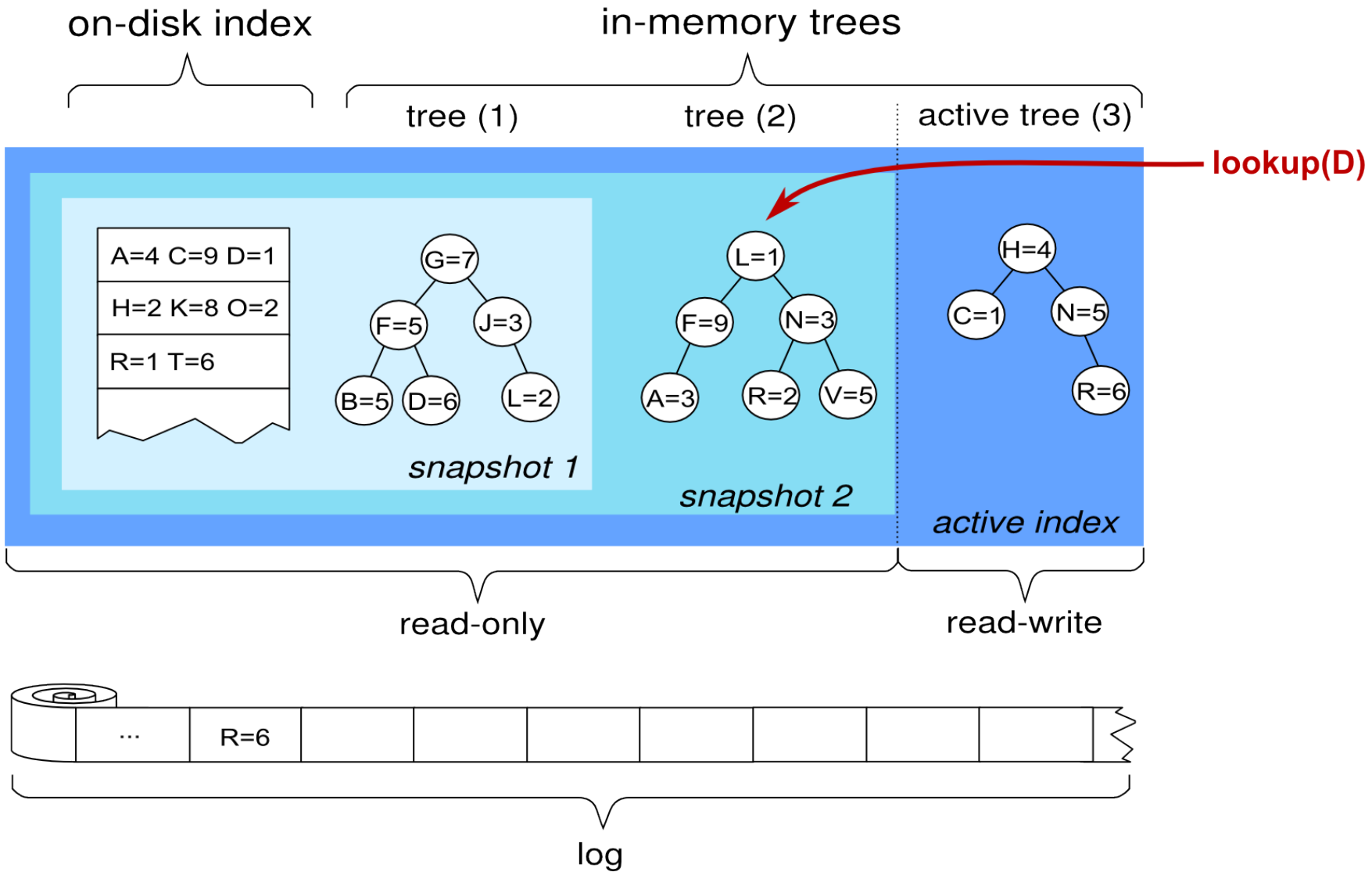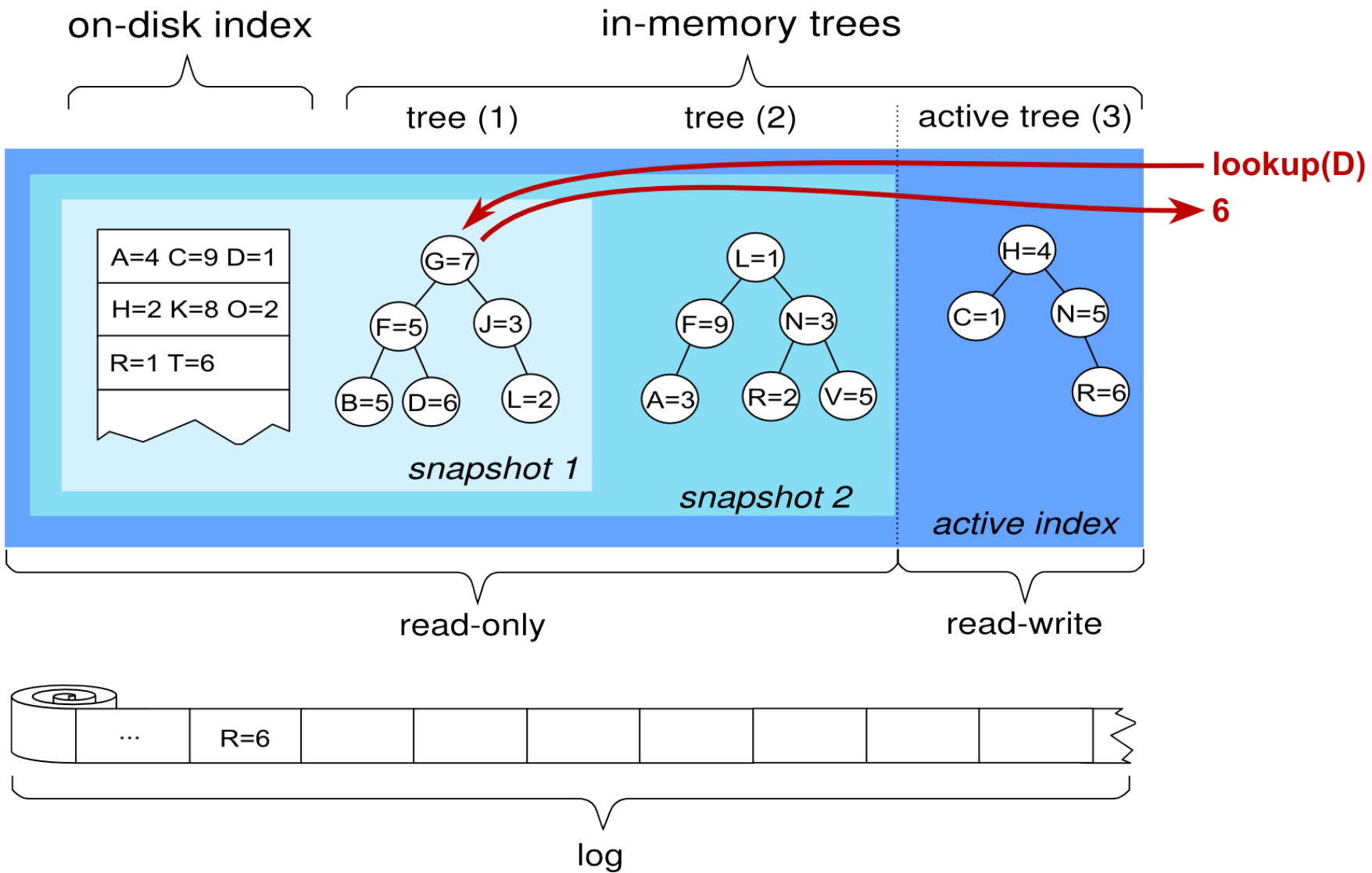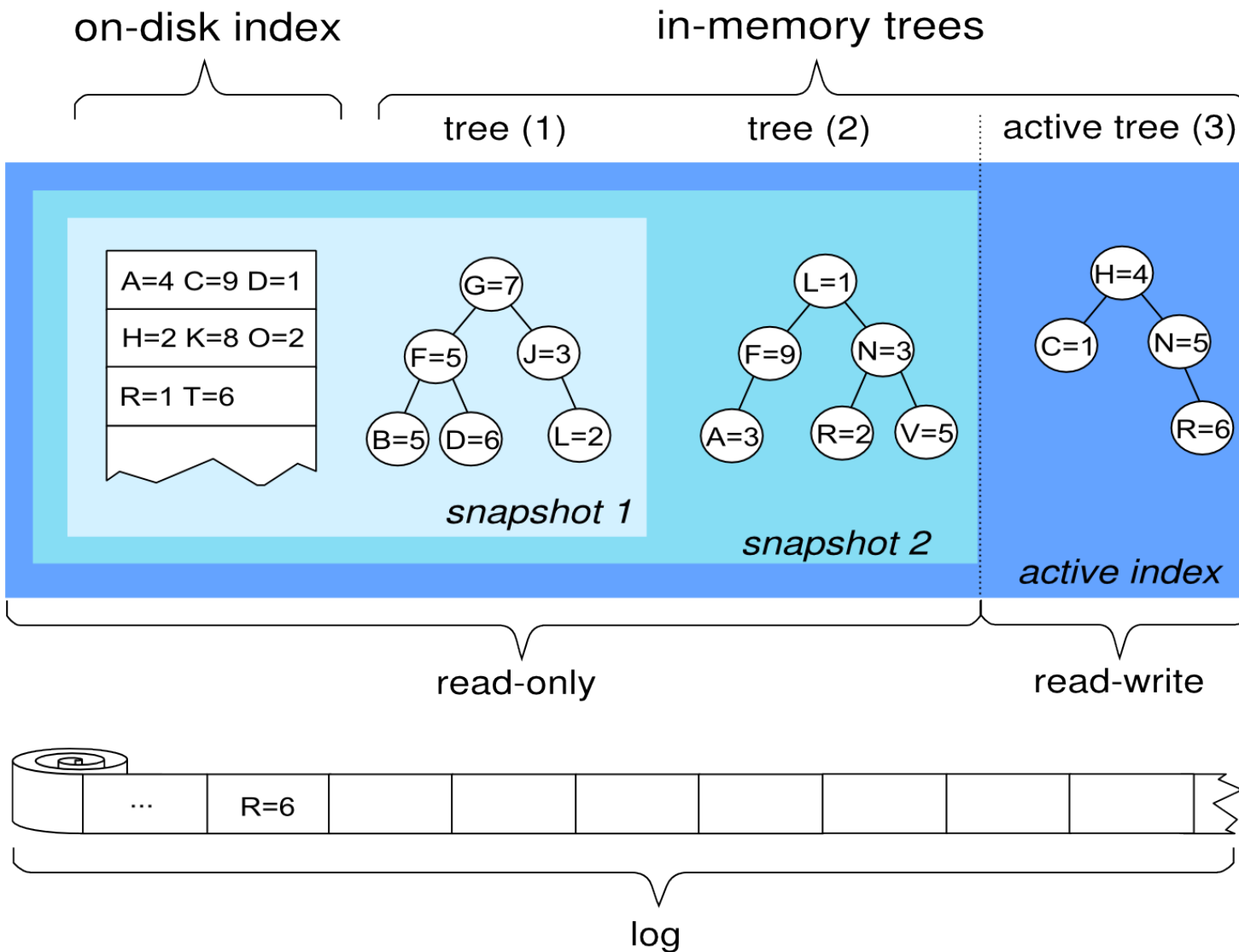
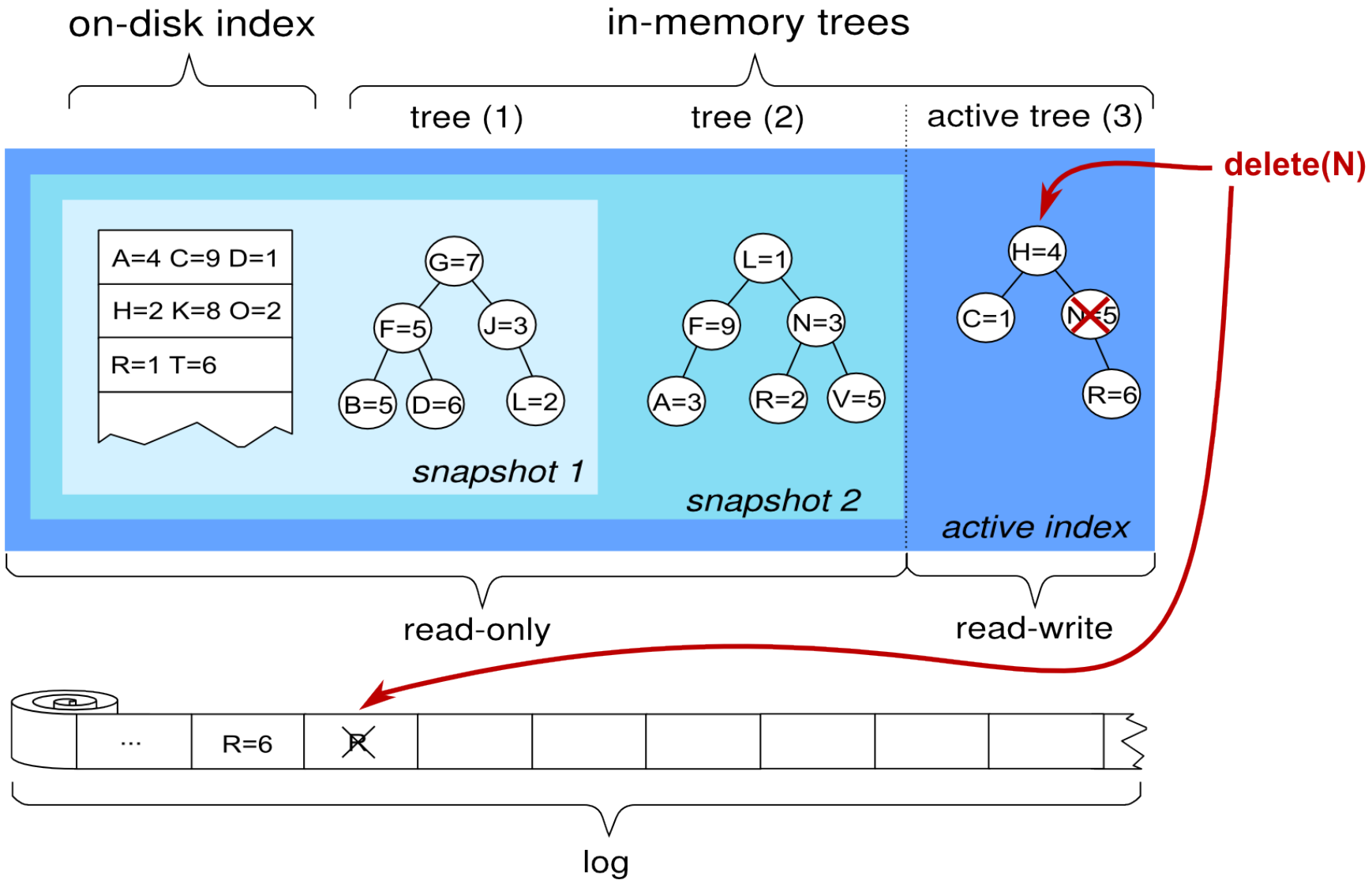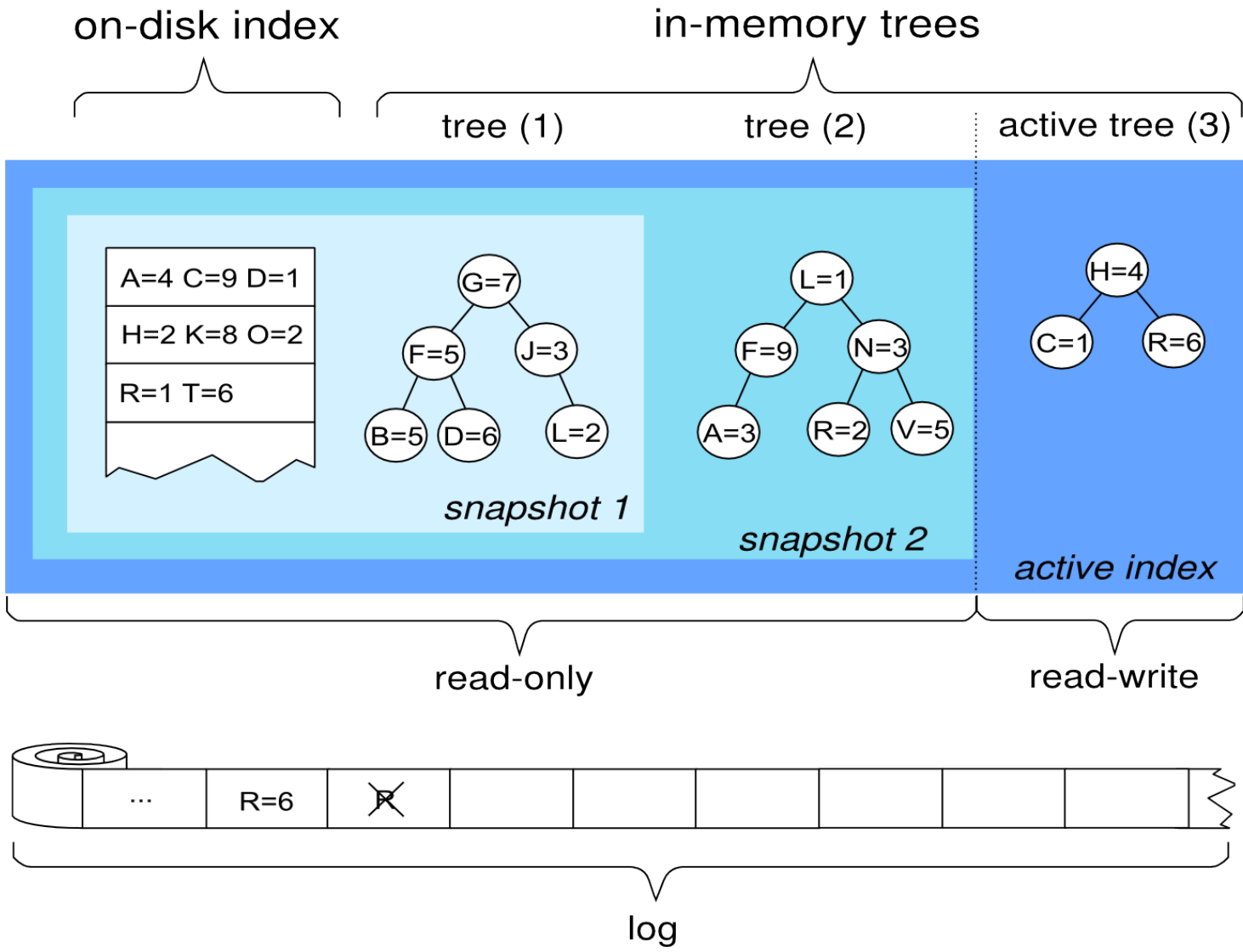# Example: Insertions

# Example: Lookups

# Example: Lookups

# Example: Lookups

# Example: Lookups

# Example: Deletions

on-disk index

in-memory trees

tree (1)          tree (2)          active tree (3)

**delete(N)**

| A=4 C=9 D=1 |
| H=2 K=8 O=2 |
| R=1 T=6 |

snapshot 1

snapshot 2

active index

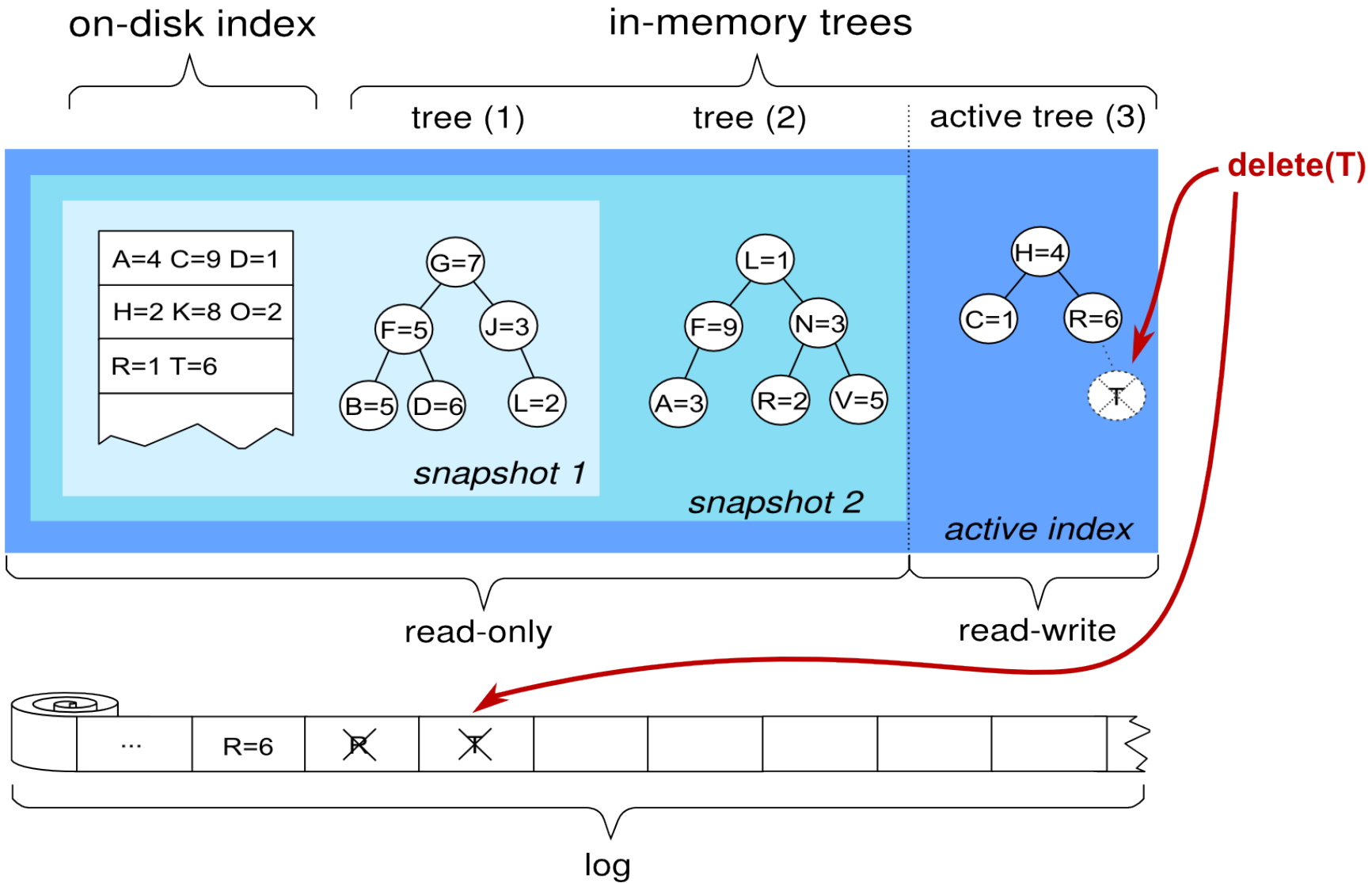read-only                    read-write

... R=6

log

# Example: Deletions

# Example: Deletions

# Example: Deletions

# Example: Range Lookups

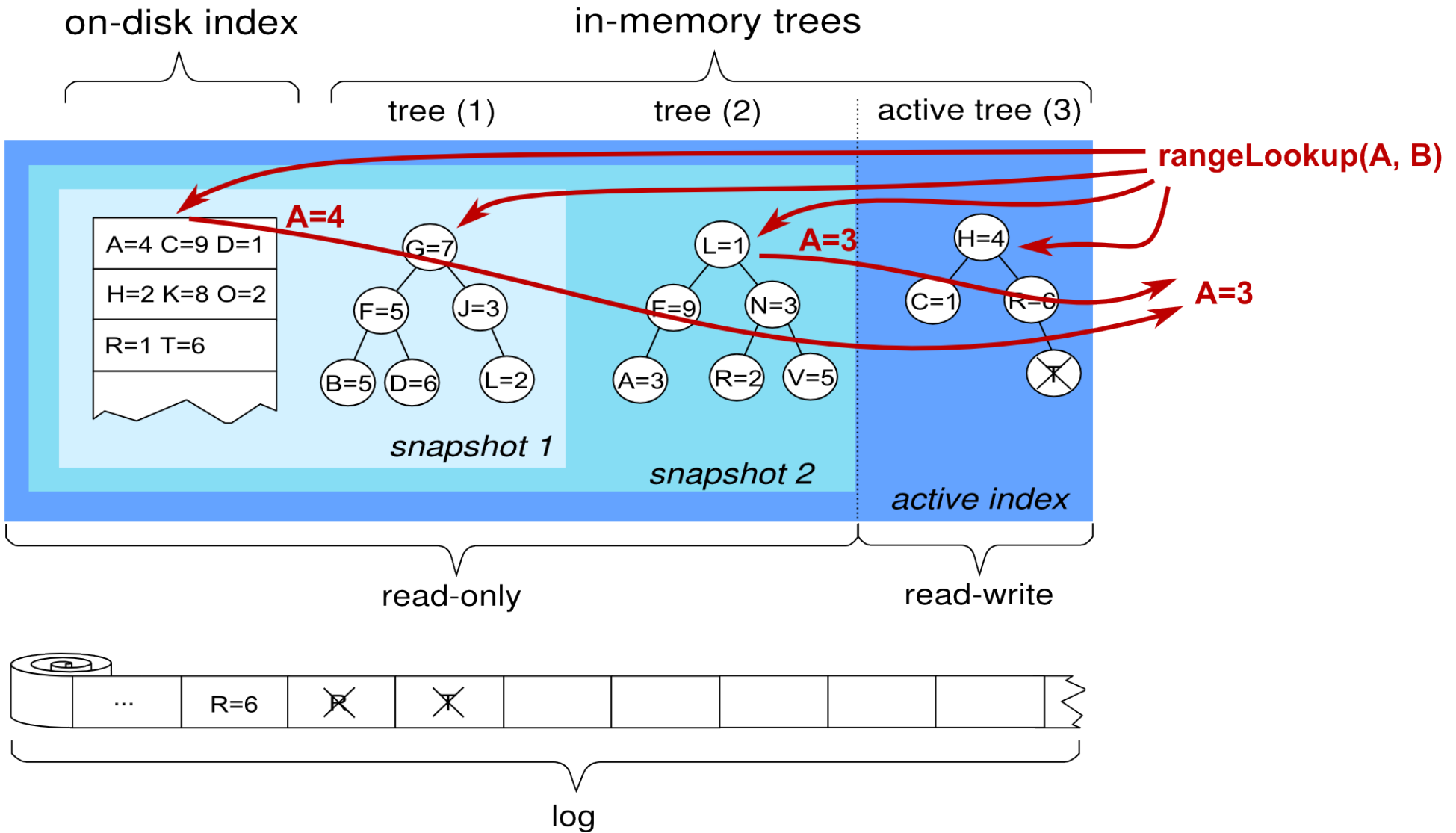# Example: Range Lookups

on-disk index

in-memory trees

tree (1)    tree (2)    active tree (3)

rangeLookup(A, B)

A=4

A=4 C=9 D=1
H=2 K=8 O=2
R=1 T=6

G=7

L=1

A=3

H=4

F=5    J=3

F=9    N=3

C=1    R=6

A=3

B=5 D=6    L=2

A=3    R=2 V=5

snapshot 1

snapshot 2

active index

read-only

read-write

... R=6

log

# Example: Range Lookups

# Example: Checkpoints

# Example: Checkpoints

# Example: Checkpoints

# Example: Checkpoints

on-disk index  in-memory trees

active tree (4)

A=3 B=5 C=1

D=6 F=9 G=7

H=4 J=3   ...

active index

read-only        read-write

log

checkpoint()

1. create snapshot
2. rangeLookup(*) at snapshot 3, write new on-disk index
3. replace on-disk index, discard snapshots, purge log

A=3 B=5 C=1

D=6 F=9 G=7

H=4 J=3   ...

- Sorted by Keys
- Block index in RAM, blocks `mmap`'ed

## BabuDB: Related Work

- Inspired by log-structured merge trees (LSM-trees)

    - Only one on-disk index

    - No „rolling merge"

- Made popular by Google Bigtable

    - Insert/lookup/merge similar as in Bigtable's Tablets

- Mapping a hierarchical directory tree to a flat database index:

# BabuDB: Advantages

– ## Why BabuDB for File System Metadata?

- Short-lived files
  - 50% of all files deleted within 5 minutes
- Atomic file system operations w/o locking or transactions
  - e.g. `rename`
- Directory content in contiguous disk regions
  - Efficient `readdir` + `stat`
- Snapshots
  - No need for multi-version data structures

– ## Linux kernel build

  – ~10M calls: 44% `stat`, 40% `open`, 15% `readlink`, 1% others

– ## Dovecot mail server + **`imaptest`**

  – ~2M calls: 51% `stat`, 48% `open`, 1% others

– ## Listing directory content

# Summary

– ## BabuDB is …

  – an efficient key-value store

  – optimized for file system metadata but also suitable for other purposes

  – suitable for large-scale databases

  – available for Java and C++ under BSD license

  – used in the XtreemFS metadata server



**http://babudb.googlecode.com**



**http://www.xtreemfs.org**

# Thank you for your attention!

# Background: XtreemFS

- **XtreemFS: a distributed replicated Internet file system**

    - part of the XtreemOS research project

    - developed since 2006 by partners from Germany, Spain and Italy



- **Object-based architecture:**

    - **MRC** stores metadata

    - **OSD**s store pure file content as objects

    - **Client**s provide POSIX file system interface

**www.xtreemfs.org**

# The XtreemOS Project

- Research project funded by the European Commission

- 19 partners from Europe and China

- XtreemFS is the data management component

  - developed by ZIB, NEC HPC Europe, Barcelona Supercomputing Center and ICAR-CNR Italy

  - ~ 3 years of development

  - first public release in August 2008

# XtreemFS: Overview

- ## What is XtreemFS?

  - a **distributed** and **replicated POSIX** compliant file system

  - **off-the-shelve** Servers – no expensive hardware

  - servers in **Java**, runs on Linux / OS X / Solaris

  - client in **C**, runs on Linux / OS X / Windows

  - **secure** (X.509 and SSL)

  - **easy** to install and maintain

  - open source (GPL)

# File System Landscape