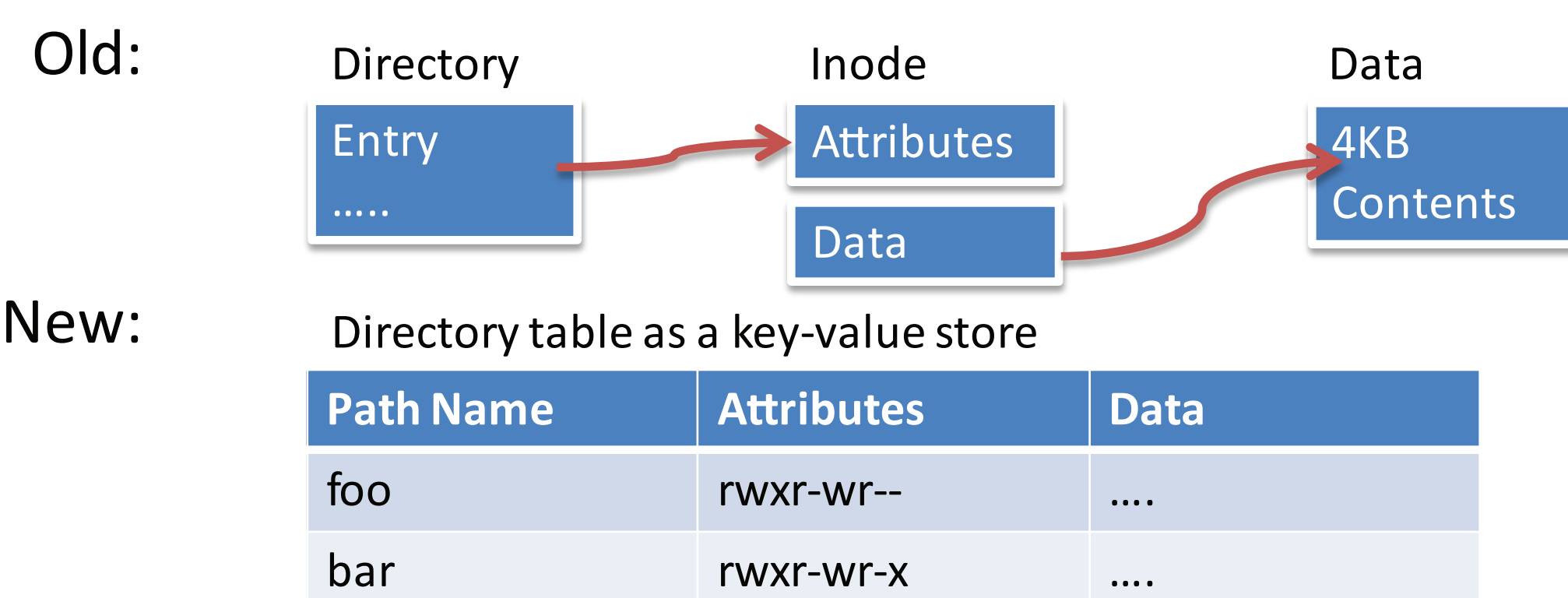
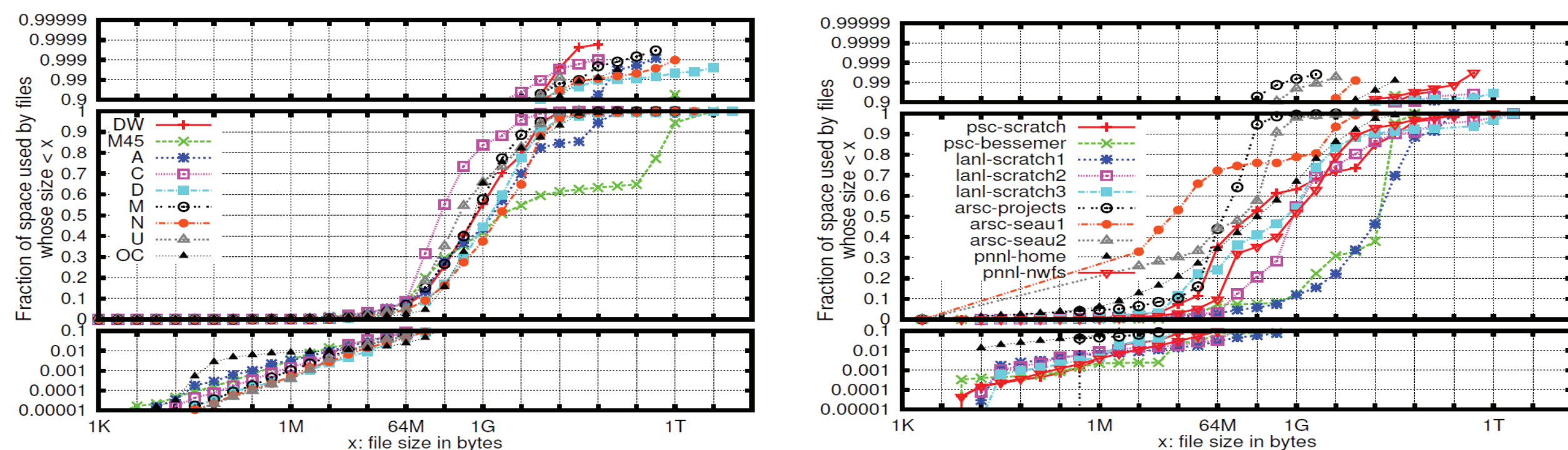


FILESYSTEM METADATA MANAGEMENT USING FAST KEY-VALUE STORE

Kai Ren, Garth Gibson (CMU)

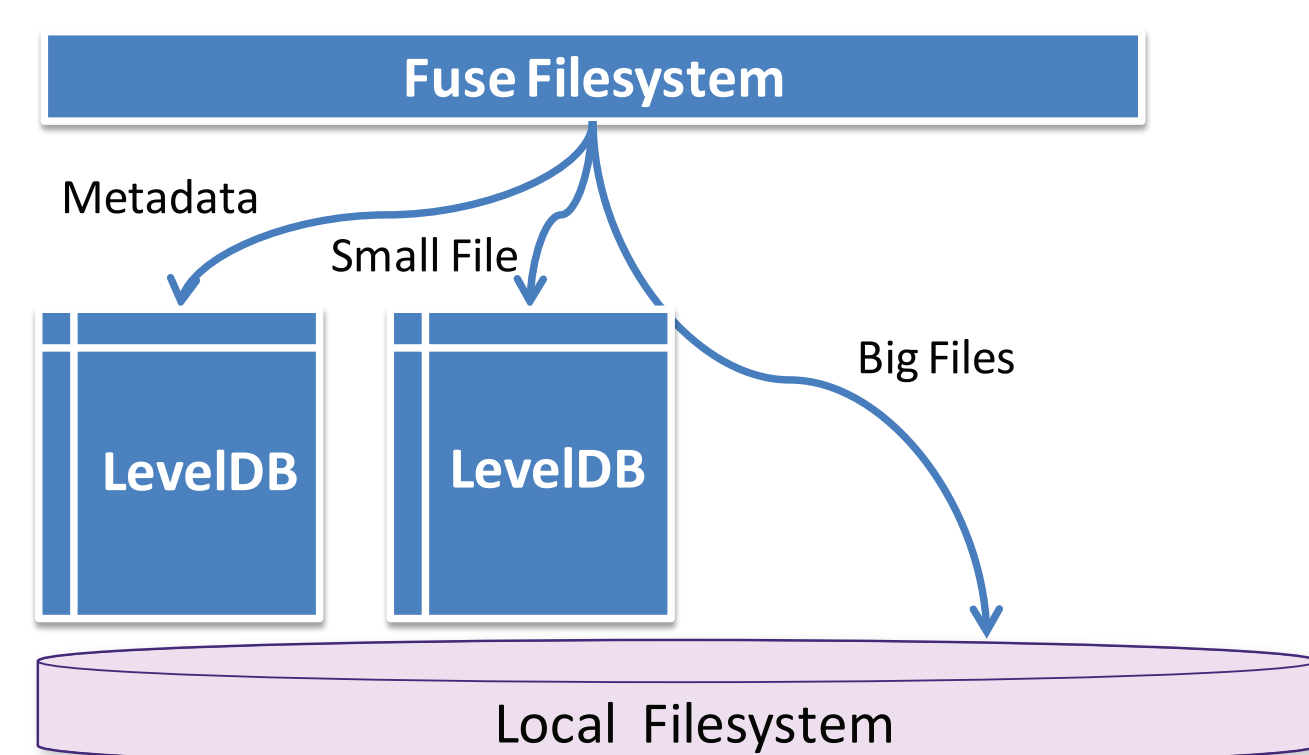
OVERVIEW

- Today's local filesystems are filled with large files but overwhelmed by more small files
- Small files widely-spread on disk needing more disk seeks to get attributes and data
- Key idea: develop new directory representation to embed file attributes and small file data
- Approach: use LevelDB with one table for all directories



FILESYSTEM DIRECTORY IMPLEMENTATION

- Built on top of Fuse file system, storing metadata and small files into LevelDB

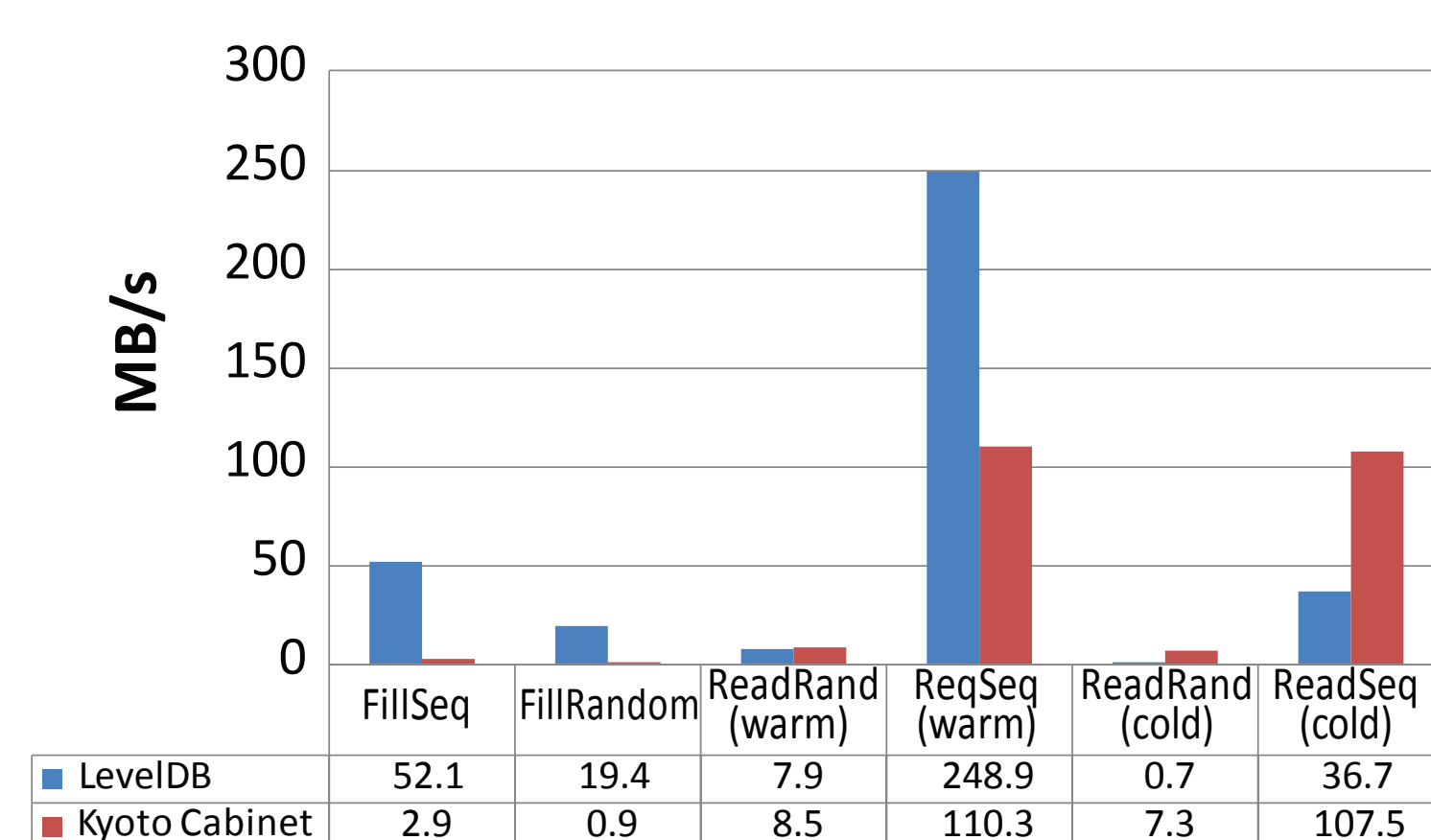


DATA LAYOUT

- Schema of Metadata:
 - Key: [Inode ID of Parent Directory, CRC Hash(filename)]
 - Value: [Filename, Attributes]
- Schema of Data:
 - Key: Inode ID, Chunk ID
 - Value: Chunk content

EVALUATION OF LEVELDB VS. B-TREES

- 1 Million entries (key: 16bytes, value: 128bytes)



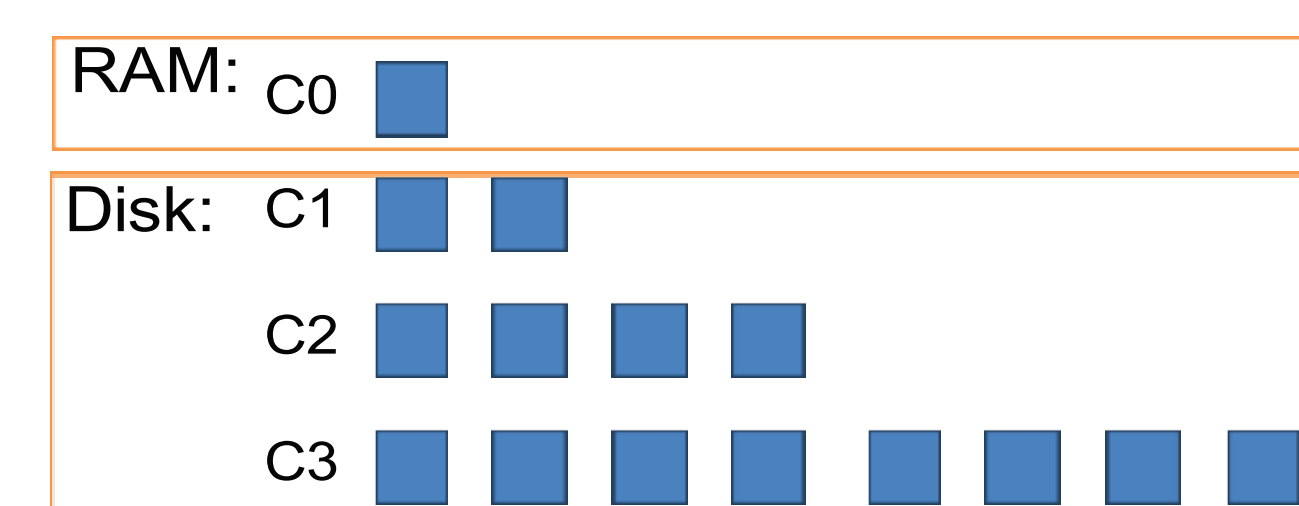
LOG-STRUCTURED MERGE TREE

- In today's Big-Data world, high ingestion workloads are prevalent and overwhelm traditional B-Trees
- Write optimized data structures like LSM (Log Structured Merge) Tree and its variants are proposed:

	Read	Update	Range Query
CoW B-Tree	$O(\log_b N)$	$O(\log_b N)$	$O(\log N + L/B)$ random I/Os
LSM Tree	$O(\log N)$	$O(\log N/B)$	$O(\log N + L/B)$ Sequential I/Os

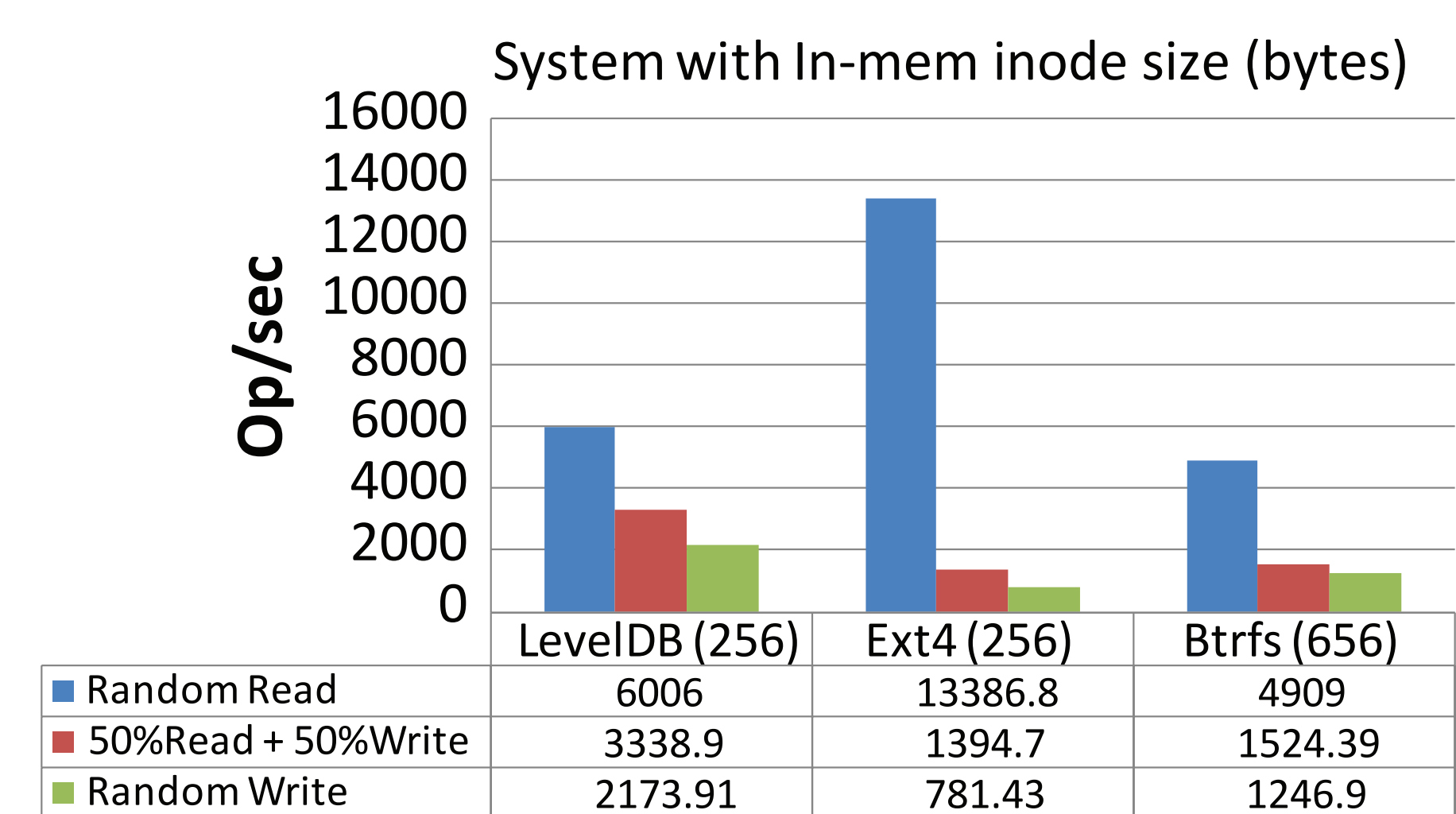
LSM DATA STRUCTURE OVERVIEW

- LSM, R-COLA, LevelDB and other variants:
 - $C[i]$ ($i > 0$): an immutable stable containing sorted (k, v) pairs, and a b-tree index
 - Size of $C[i+1] = R * \text{size of } C[i]$
- Insertion: write to in-ram tablet, later will be written to disk
- Compaction: merging $C[0] \dots C[i]$ into $C[i+1]$ when $C[0] \dots C[i]$ are full
- Range query: search sstables with appropriate range in each $C[i]$



EVALUATION OF FILESYSTEM METADATA

- Workload: 1M empty files, only attributes in LevelDB
- Random Read: stat a file randomly
- Random Write: chmod / utime a file randomly
- Memory limit: 600MB for LevelDB, 100MB for Ext4 and Btrfs
 - Limiting memory tries to control different internal cache sizes.



- Workload: Write 1M 1KB files sequential, overwrite 1M 1KB files. Compare sequential overwrite, random overwrite, and overwrite with read operations
- Memory Limit: 512MB, Disk Partition: 5GB

