

DBLK: Deduplication for primary storage

Yoshi Tsuchiya and Takashi Watanabe
Fujitsu

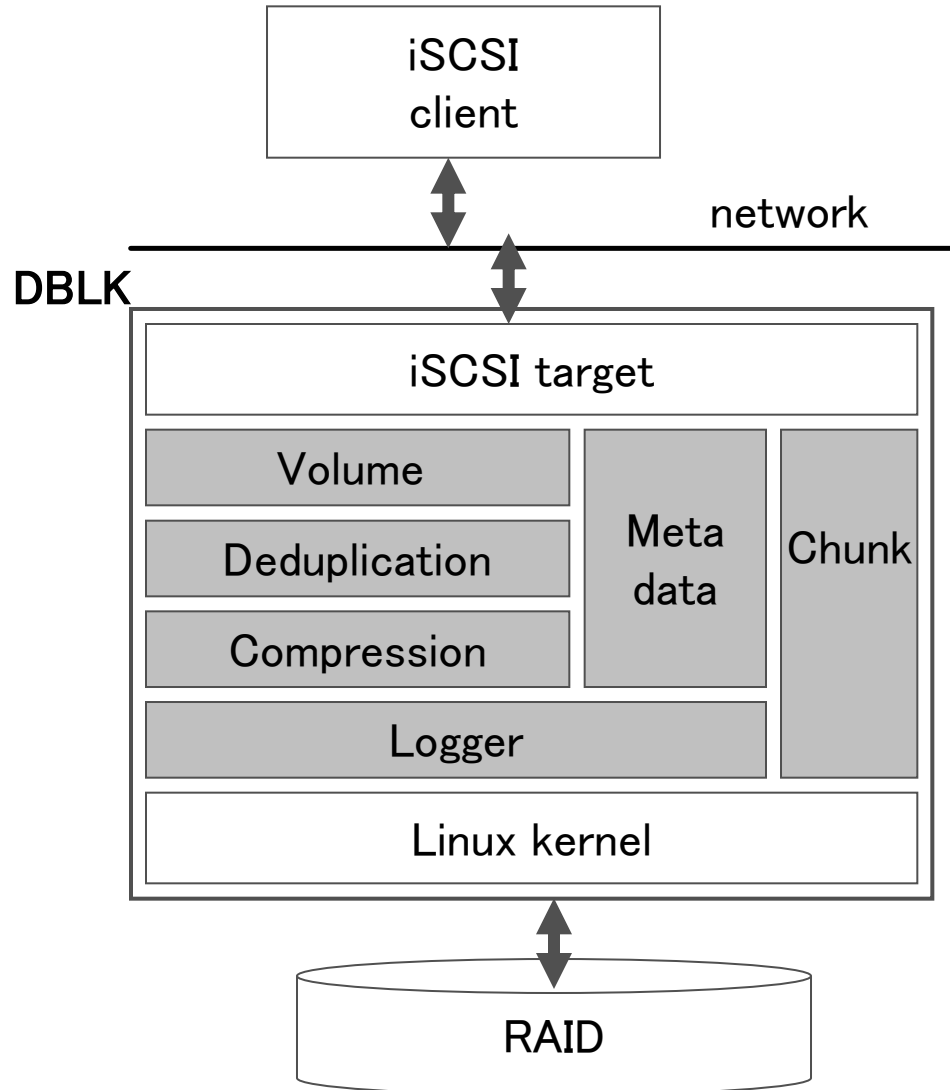
5/26/2011 IEEE MSST Research Track

This work is supported by the Green IT project of New Energy and Industrial Technology Development Organization (NEDO)

- Background
- DBLK outline
- DBLK Hash index
- Multilayer Bloom Filter (MBF)
- MBF Bitwise Transposition (MBF BT)
- Summary

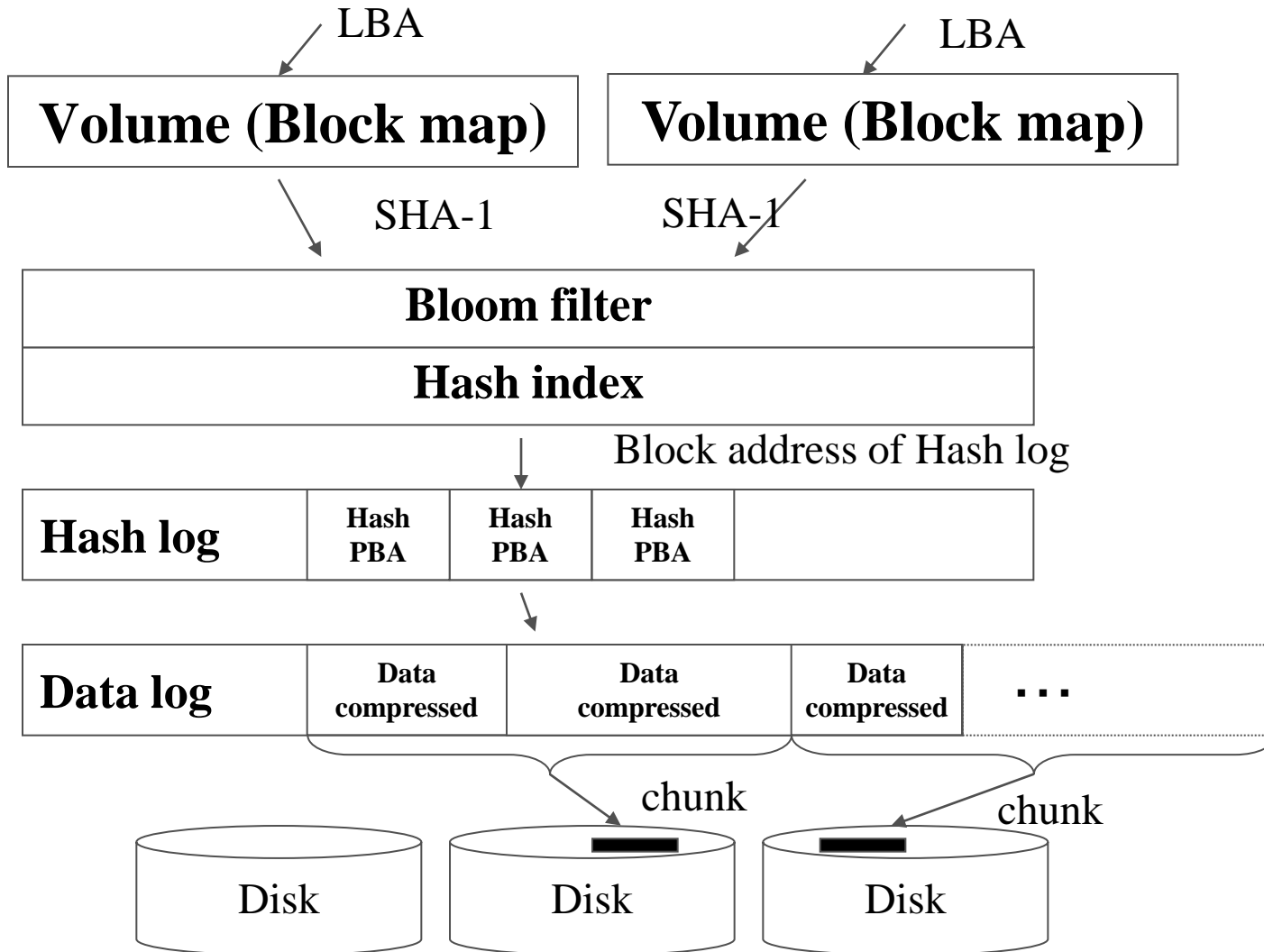
- Deduplication: Reduce size of storage by removing redundancy in data
- Dedup for Primary storage/Backup storage
 - Primary storage: file servers, small latency is required
- Inline/Post-process dedup
 - Inline: immediately deduplicate and write data
 - Post process: write data first, deduplicate later, needs extra storage
- Deduplication unit
 - Fixed size block
 - File
 - Variable-sized segment

DBLK iSCSI implementation

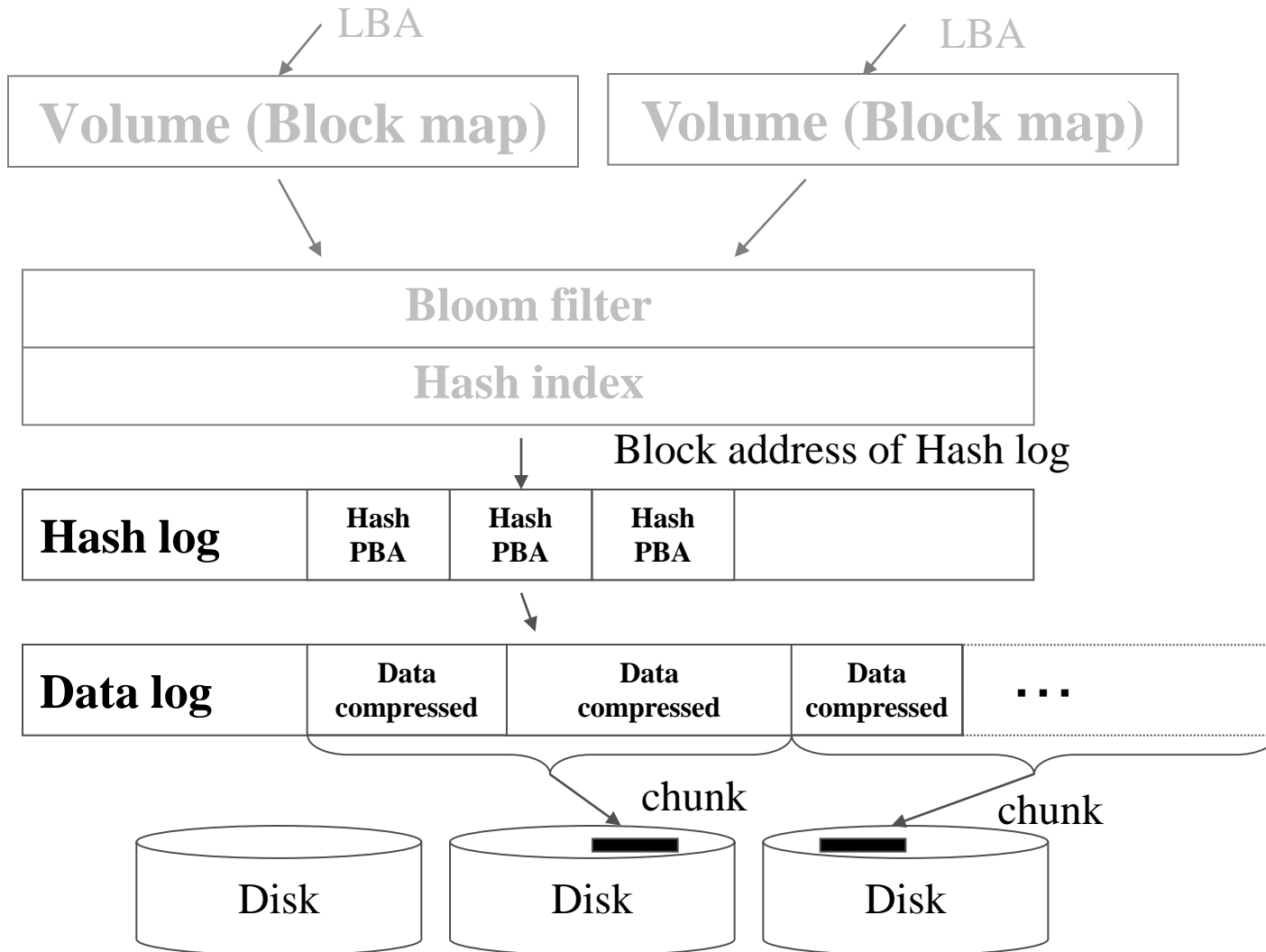


**Primary block storage
Inline, block-wise deduplication,
thin provisioning and compression**

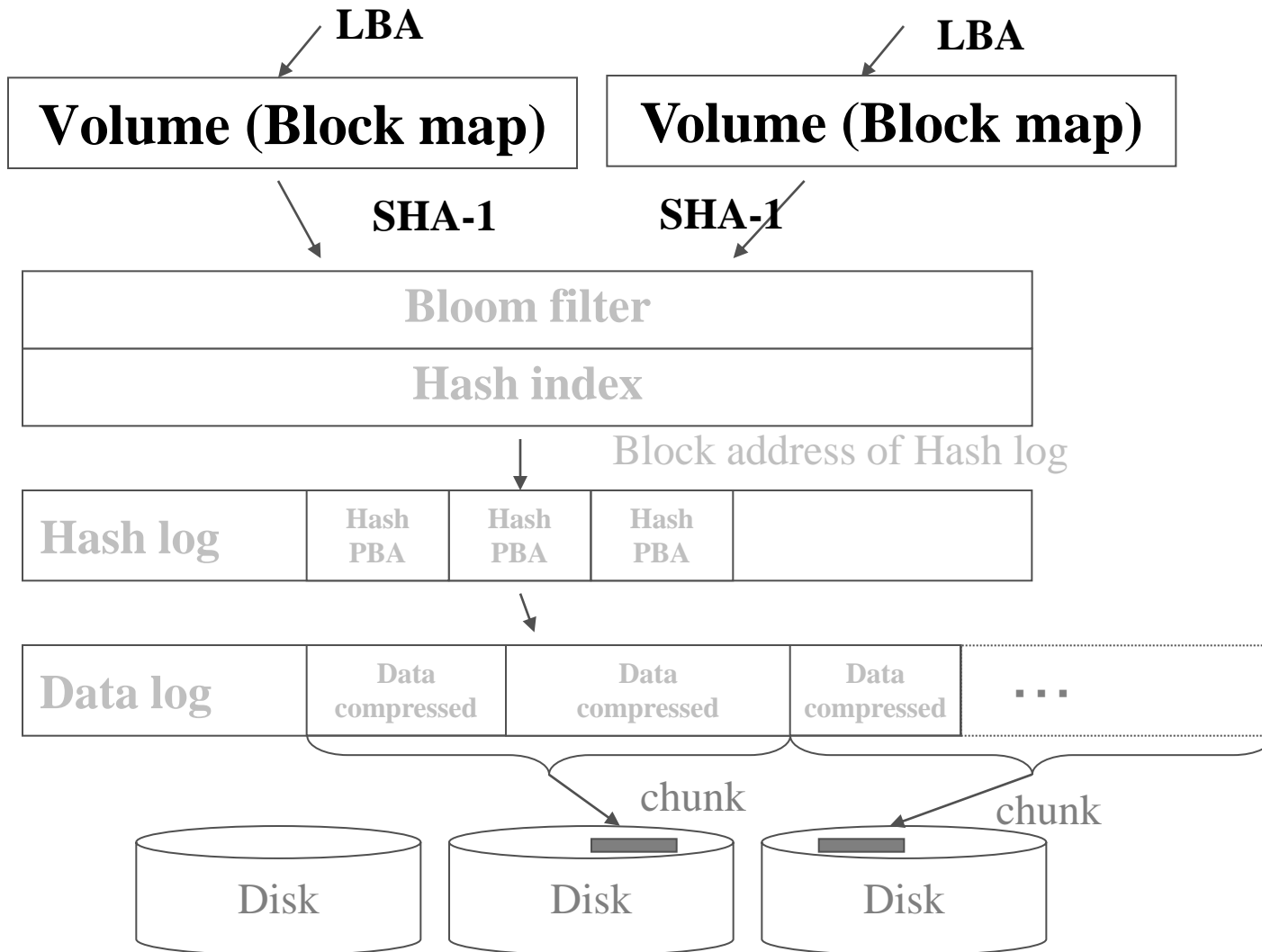
DBLK Data Structure



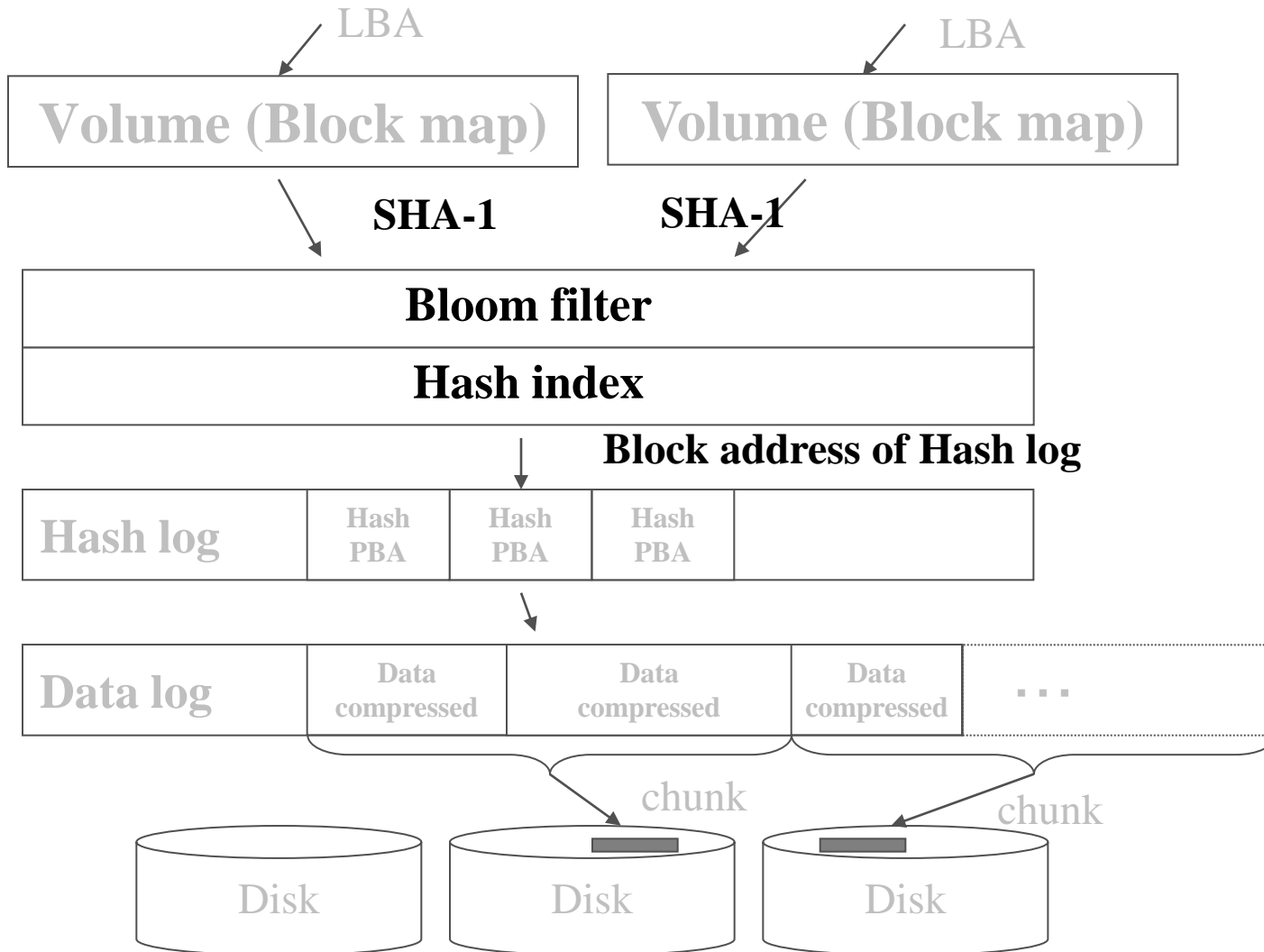
DBLK Data Structure: Logs



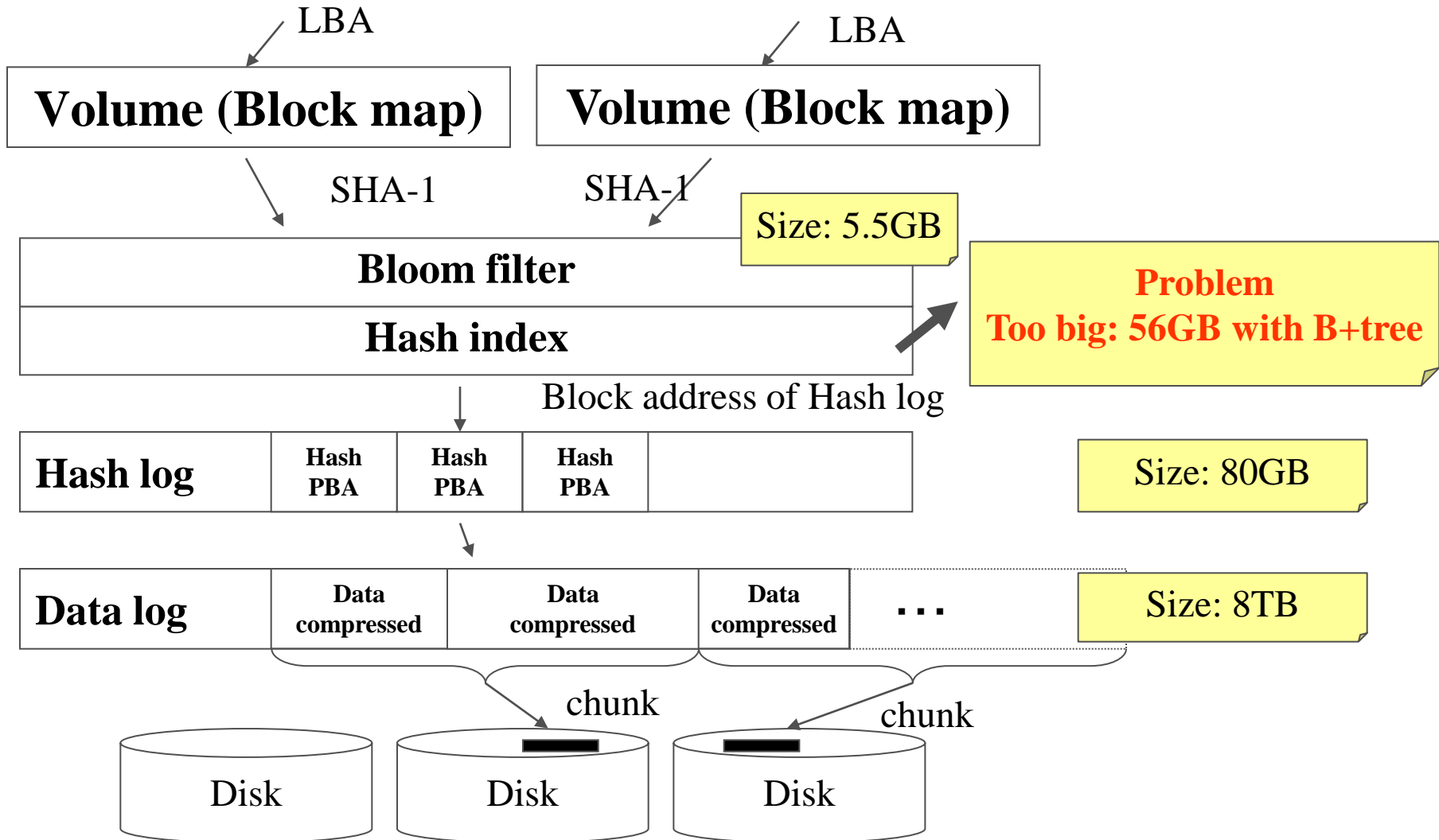
DBLK Data Structure: Block maps



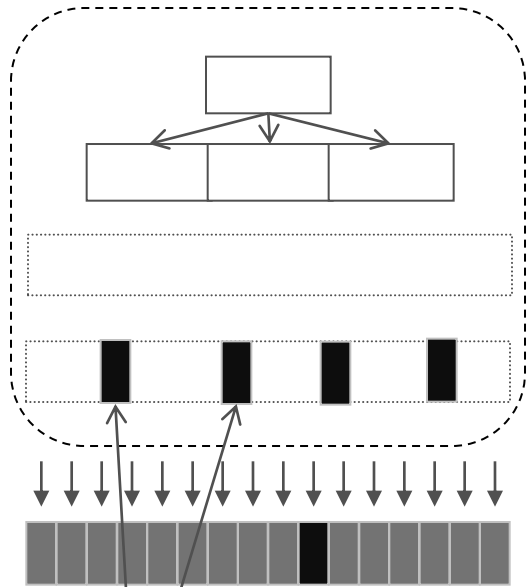
DBLK Data Structure: Hash index



DBLK Data Structure



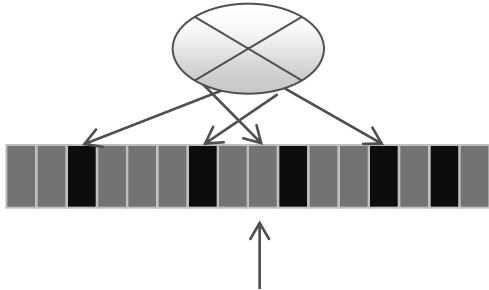
B+tree



I/Os for leaf nodes

Append only log.
only one dirty block

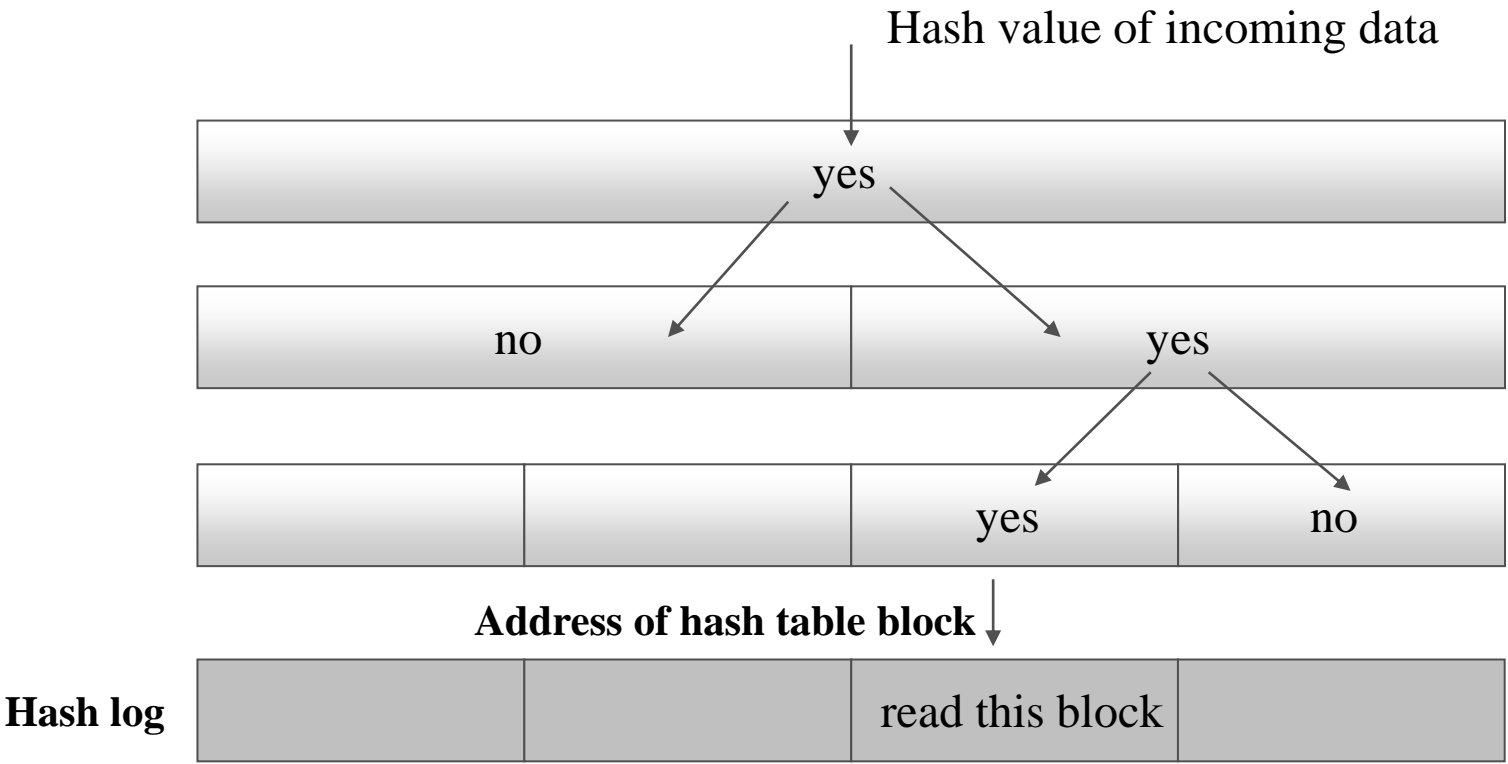
Hash function

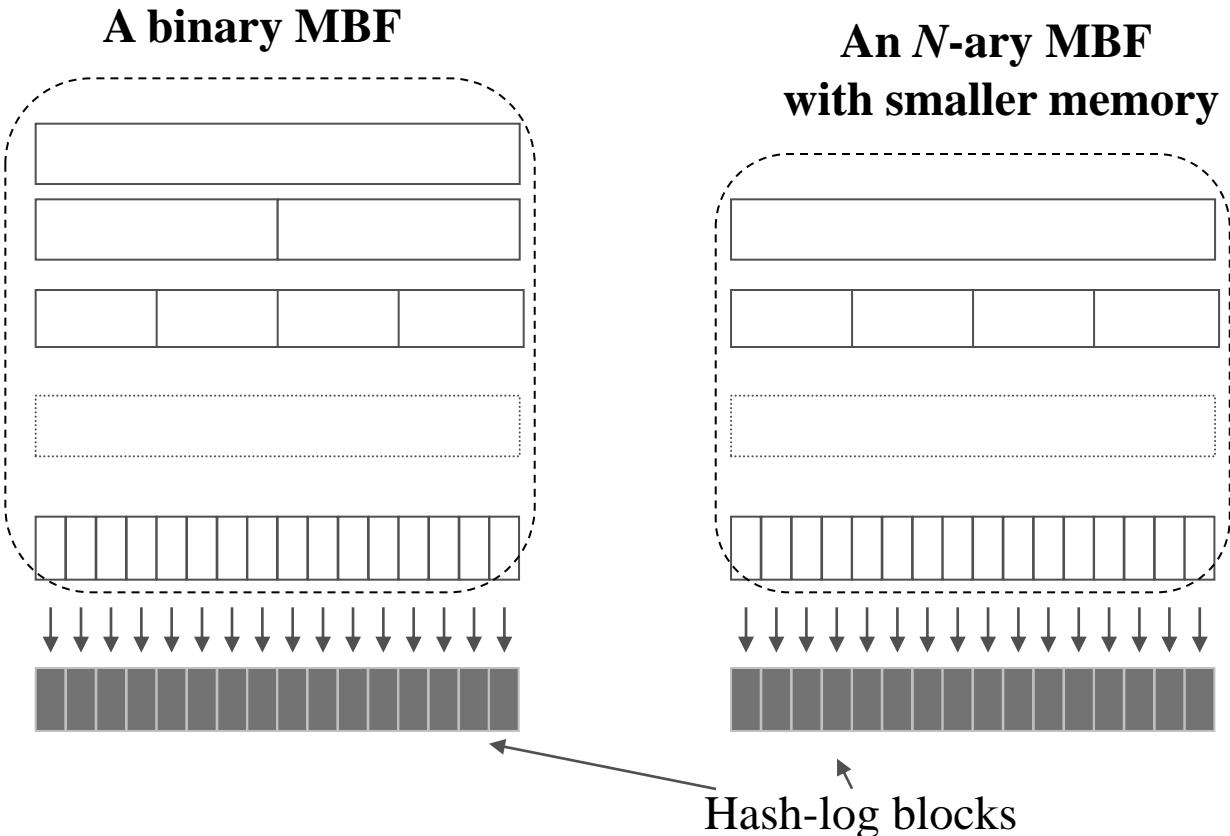


Hash bucket with
lots of I/O

Solution: Multilayer Bloom filter

- Layers of Bloom filters
 - If it reaches to a lowest filter, read the hash-log block
- False positive means
 - Extra I/Os but a correct answer

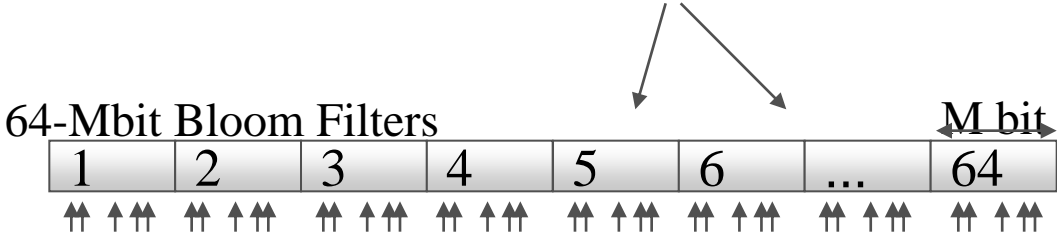




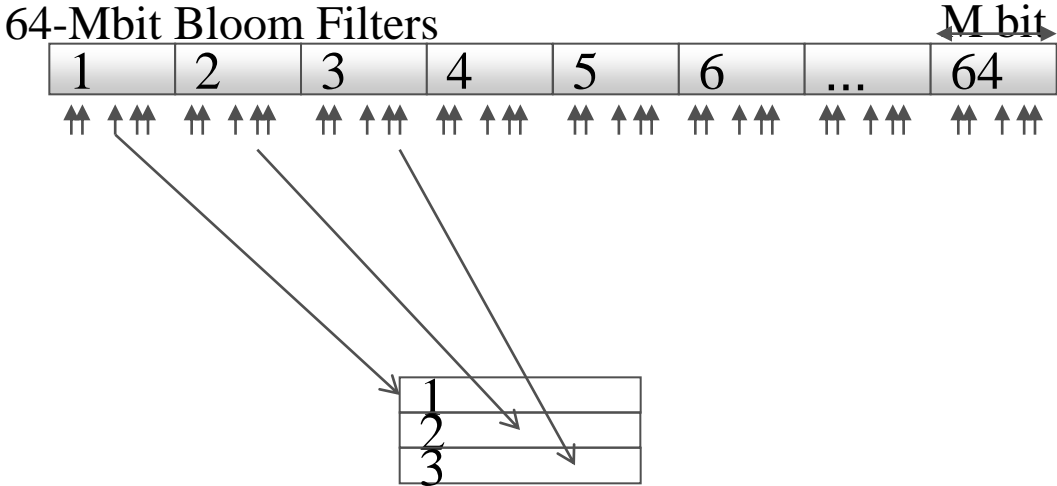
The bigger N , the shorter MBF = less memory.

MBF Optimization: MBF-BT

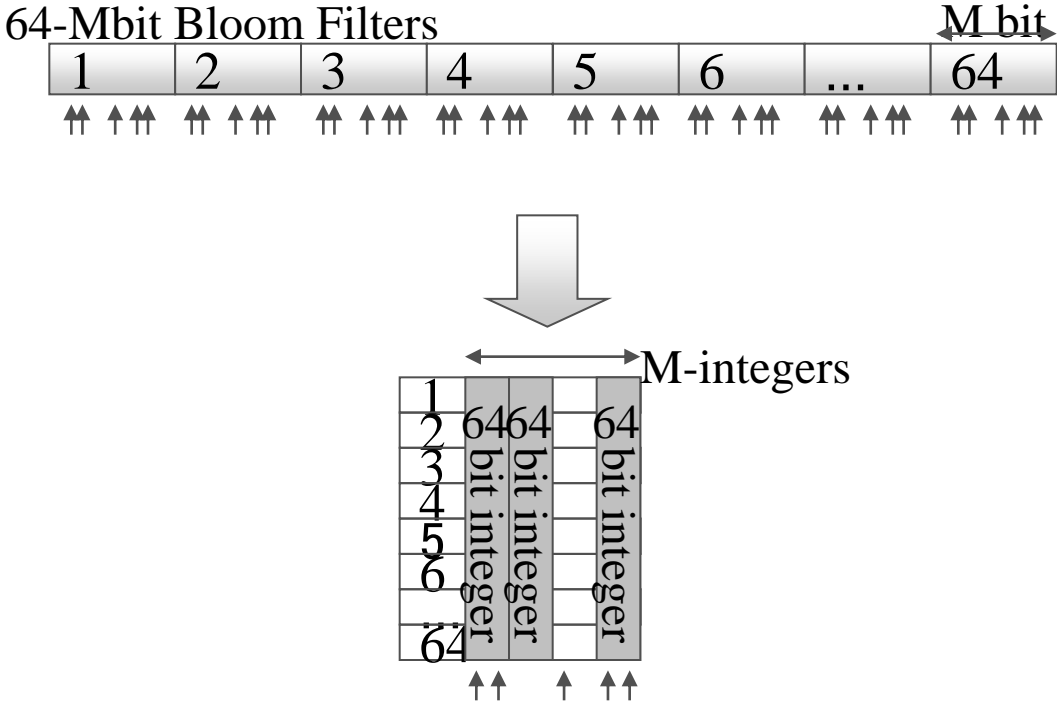
Compute (b1&b2&...&b10) 64 times



MBF Optimization: MBF-BT

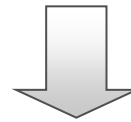
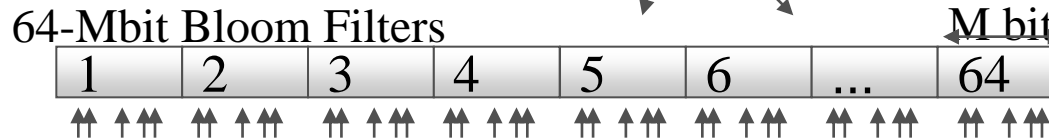


MBF Optimization: MBF-BT

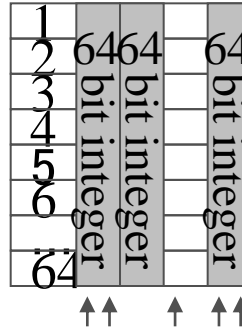


MBF Bitwise Transposition

Compute (b1&b2&...&b10) 64 times



M-integers

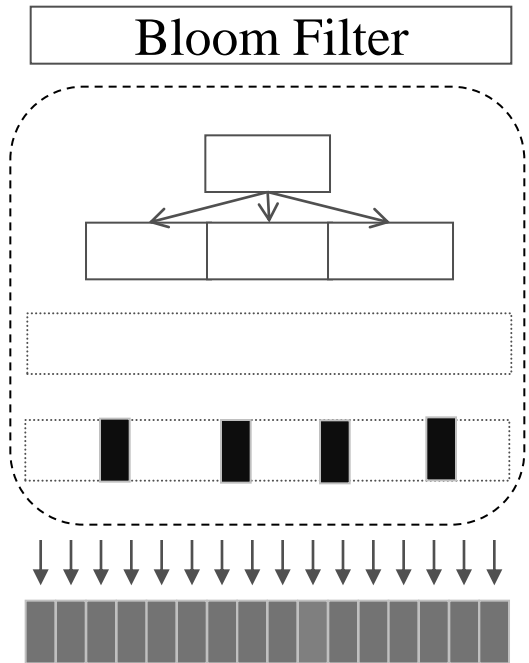


Compute $x_1 \& x_2 \& \dots \& x_{10}$ once
64x speedup

Bitwise transpose to
M-64bit integers

Large N for N -ary MBF
1664 in our implementation

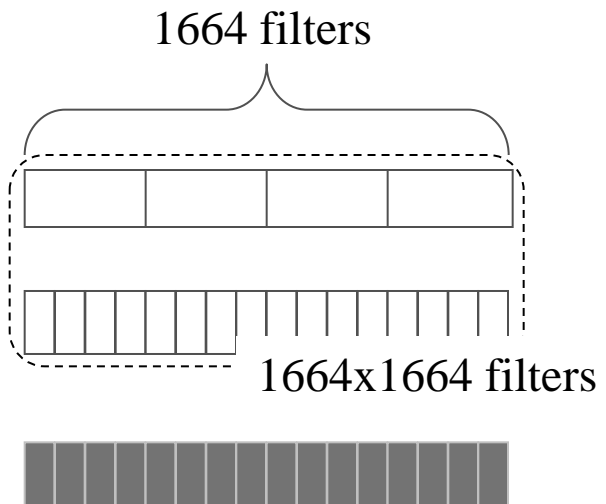
Bloom filter and B+tree



5.5GB+56GB



Result: 2-layer MBF



5.5GB for each layer, i.e. 11GB

■ MBF

- On-memory data structure
- Helps reduce the latency of primary dedup


■ MBF bitwise transposition

- Reduces CPU cost to lookup bunch of Bloom filters
- Optimizes access to MBF

■ DBLK

- Smaller latency than other implementation
- No I/O for hash index

Thank you!



FUJITSU

shaping tomorrow with you