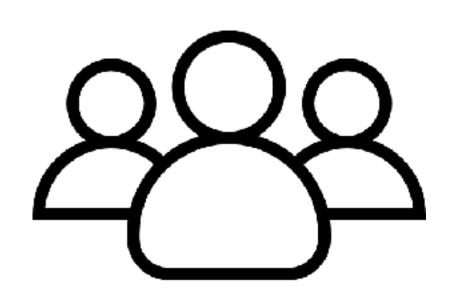
Tutorial on LSM-trees

Database System Technology

Groups

Form for registration was sent today:)

Complete by Sunday





Question 1 - picking a storage engine

We would like to choose KV-store engine between BerkeleyDB vs. RocksDB for a workload consisting of 10% random gets and 90% random puts. BerkeleyDB uses a B-tree. RocksDB uses a leveled LSM-tree with size ratio T=10 and a buffer size of P=2²⁶ entries. Assume N=2⁴⁰ and B=2⁷. All internal nodes fit in memory. We are using a disk drive. What is your choice and why?







Question 1 - picking a storage engine

We would like to choose KV-store engine between BerkeleyDB vs. RocksDB for a workload consisting of 10% random gets and 90% random puts. BerkeleyDB uses a B-tree. RocksDB uses a leveled LSM-tree with size ratio T=10 and a buffer size of P=2²⁶ entries. Assume N=2⁴⁰ and B=2⁷. All internal nodes fit in memory. We are using a disk drive. What is your choice and why?

B-tree: Each "put" costs 1 read & 1 write I/O. Each get costs 1 read I/O.

As we are using disk, read/write I/O costs are symmetric.

Avg. #I/O per operation: $0.9 * 2 + 0.1 * 1 \approx 1.9$



Memory

Storage

Memory

Memor

Question 1 - picking a storage engine

We would like to choose KV-store engine between BerkeleyDB vs. RocksDB for a workload consisting of 10% random gets and 90% random puts. BerkeleyDB uses a B-tree. RocksDB uses a leveled LSM-tree with size ratio T=10 and a buffer size of P=2²⁶ entries. Assume N=2⁴⁰ and B=2⁷. All internal nodes fit in memory. We are using a disk drive. What is your choice and why?

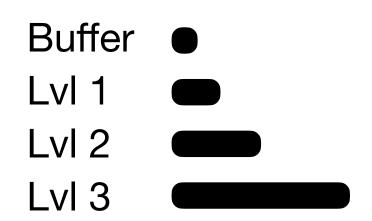
B-tree: Avg. #I/O per operation: $0.9 * 2 + 0.1 * 1 \approx 1.9$

LSM-tree: A put costs 2 * (T/B) · log_T(N/P) ≈ 0.66 read & write I/Os

A get costs $log_T(N/P) \approx 4.2 \text{ read } I/Os$

Avg. # I/Os per operation: $0.66 \cdot 0.9 + 4.2 \cdot 0.1 \approx 1.0 \text{ I/Os}$ (Cheaper)





A friend tells you they switched from using a B-tree to a basic LSM-tree (with size ratio 2), yet write-amplification actually increased. There are N=2⁴⁰ entries, and the memtable size P and block size B are both B=P=2⁵ entries. Explain why this happened. Identify three tuning options for reducing write-amplification and the trade-off of each one of them.

Buffer Lvl 1 Lvl 2 Lvl 3

A friend tells you they switched from using a B-tree to a basic LSM-tree (with size ratio 2), yet write-amplification actually increased. There are $N=2^{40}$ entries, and the memtable size P and block size B are both $B=P=2^5$ entries. Explain why this happened. Identify three tuning options for reducing write-amplification and the trade-off of each one of them.

LSM-tree before tuning: $(1/B) * log_2(N/P) = 1.1$ write I/O (Costlier than B-tree)

Buffer Lvl 1 Lvl 2 Lvl 3

A friend tells you they switched from using a B-tree to a basic LSM-tree (with size ratio 2), yet write-amplification actually increased. There are N=2⁴⁰ entries, and the memtable size P and block size B are both B=P=2⁵ entries. Explain why this happened. Identify three tuning options for reducing write-amplification and the trade-off of each one of them.

LSM-tree before tuning: $(1/B) * log_2(N/P) = 1.1$ write I/O (Costlier than B-tree)

Increasing the buffer size P reduces the number of levels and thus WA.



A friend tells you they switched from using a B-tree to a basic LSM-tree (with size ratio 2), yet write-amplification actually increased. There are $N=2^{40}$ entries, and the memtable size P and block size B are both $B=P=2^5$ entries. Explain why this happened. Identify three tuning options for reducing write-amplification and the trade-off of each one of them.

LSM-tree before tuning: $(1/B) * log_2(N/P) = 1.1$ write I/O (Costlier than B-tree)

Increasing the buffer size P reduces the number of levels and thus WA.



We can increase the buffer size to decrease the number of levels across which entries get merged. E.g., $P=2^{20}$ entries leads to:

$$(1/B) * log_2(N/P) = 0.63$$

A friend tells you they switched from using a B-tree to a basic LSM-tree (with size ratio 2), yet write-amplification actually increased. There are $N=2^{40}$ entries, and the memtable size P and block size B are both $B=P=2^5$ entries. Explain why this happened. Identify three tuning options for reducing write-amplification and the trade-off of each one of them.

LSM-tree before tuning: $(1/B) * log_2(N/P) = 1.1$ write I/O (Costlier than B-tree)

Increasing the buffer size P reduces the number of levels and thus WA.



We can increase the buffer size to decrease the number of levels across which entries get merged. E.g., $P=2^{20}$ entries leads to:

 $(1/B) * log_2(N/P) = 0.63$

Trade-off: requires more memory

A friend tells you they switched from using a B-tree to a basic LSM-tree (with size ratio 2), yet write-amplification actually increased. There are N=2⁴⁰ entries, and the memtable size P and block size B are both B=P=2⁵ entries. Explain why this happened. Identify three tuning options for reducing write-amplification and the trade-off of each one of them.

We can further employ tiering, say, with size ratio T=10 e.g., $(1/B) * log_T(N/P) = (1/2^5) * log_{10}(2^{40}/2^{20}) = 0.19$

Lvl 2

A friend tells you they switched from using a B-tree to a basic LSM-tree (with size ratio 2), yet write-amplification actually increased. There are $N=2^{40}$ entries, and the memtable size P and block size B are both $B=P=2^5$ entries. Explain why this happened. Identify three tuning options for reducing write-amplification and the trade-off of each one of them.

We can further employ tiering, say, with size ratio T=10

e.g.,
$$(1/B) * log_{\tau}(N/P) = (1/2^5) * log_{\tau}(2^{40}/2^{20}) = 0.19$$



A friend tells you they switched from using a B-tree to a basic LSM-tree (with size ratio 2), yet write-amplification actually increased. There are N=2⁴⁰ entries, and the memtable size P and block size B are both B=P=2⁵ entries. Explain why this happened. Identify three tuning options for reducing write-amplification and the trade-off of each one of them.

We can further employ tiering, say, with size ratio T=10

e.g.,
$$(1/B) * log_{\tau}(N/P) = (1/2^5) * log_{10}(2^{40}/2^{20}) = 0.19$$

Buffer Trade-off: reads get more expensive.

Lvl 2

Last approach: Make the LSM-tree unclustered. More in next question.



In a clustered LSM-trees, the values are stored within the LSM-tree alongside their keys. Another approach is an unclustered LSM-tree, which stores values in an append-only file and indexes them using key-pointer pairs from within the LSM-tree. What's the impact of clustered vs. unclustered LSM-trees on put/get/scan performance. Propose adjusted cost models. Assume a basic LSM-tree (size ratio T=2), insertions only (no deletes or updates), and a 1 page memtable. Let K be the number of key-pointer pairs fitting into a page, and let B be the number of key-value pairs fitting in a page. Based on your models, identify cases where each of the two approaches shines.

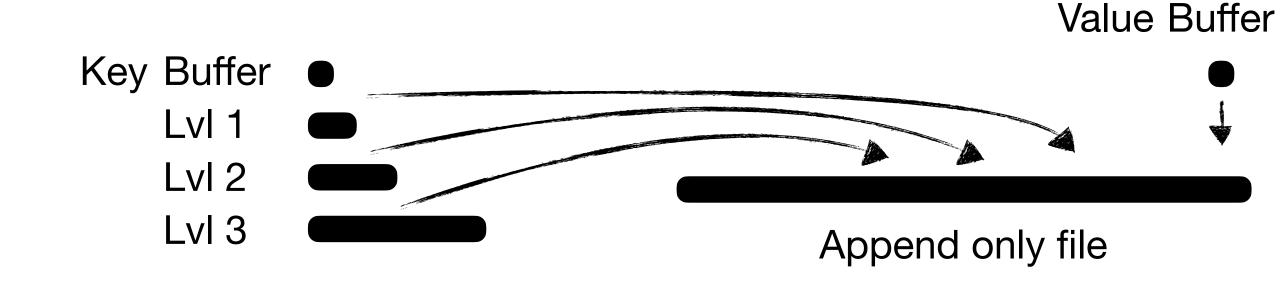
Clustered	Unclustered?
Olugici ca	Oliciusieleu:

#levels: $L_C = log_T(N/B)$

put: $O(L_C / B)$

get: O(L_C)

scan: $O(L_C+S/B)$



	Clustered	Unclustered?		
#levels:	$L_C = log_T(N/B)$	L _U =log _T (N/K)	Key Buffer ●	Value Buffer
put:	O(L _C / B)		Lvl 2	
get:	O(L _C)		Lvl 3	Append only file
scan:	O(L _C +S/B)			

	Clustered	Unclustered?		
#levels:	L _C =log _T (N/B)	Lu=log _T (N/K)	Key Buffer ■	Value Buffer
put:	O(L _C / B)	$O(1/B + L_U/K)$	Lvl 1	
get:	O(L _C)		Lvl 3	Append only file
scan:	O(L _C +S/B)			

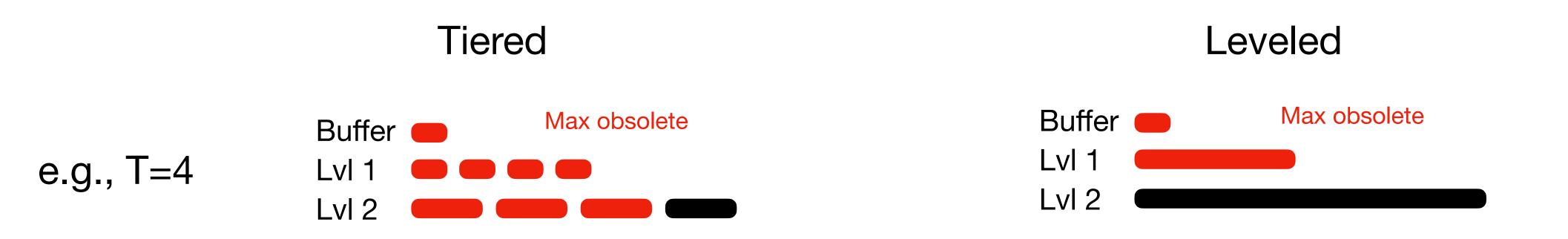
	Clustered	Unclustered?		
#levels:	L _C =log _T (N/B)	Lu=log _T (N/K)	Key Buffer ■	Value Buffer
put:	O(L _C / B)	$O(1/B + L_U/K)$	Lvl 2	
get:	O(L _C)	O(L _U + 1)	Lvl 3	Append only file
scan:	O(L _C +S/B)			

	Clustered	Unclustered		
#levels:	L _C =log _T (N/B)	Lu=log _T (N/K)	Key Buffer ■	Value Buffer
put:	O(L _C / B)	$O(1/B + L_U/K)$	Lvl 1 Lvl 2	
get:	O(L _C)	O(L _U + 1)	Lvl 3	Append only file
scan:	O(L _C +S/B)	O(L _U +S)		

	Clustered	Unclustered	Unclustered is better when
#levels:	L _C =log _T (N/B)	Lu=log _T (N/K)	
put:	O(L _C / B)	O(1/B + L _U /K)	
get:	O(L _C)	O(L _U + 1)	
scan:	O(L _C +S/B)	O(L _U +S)	

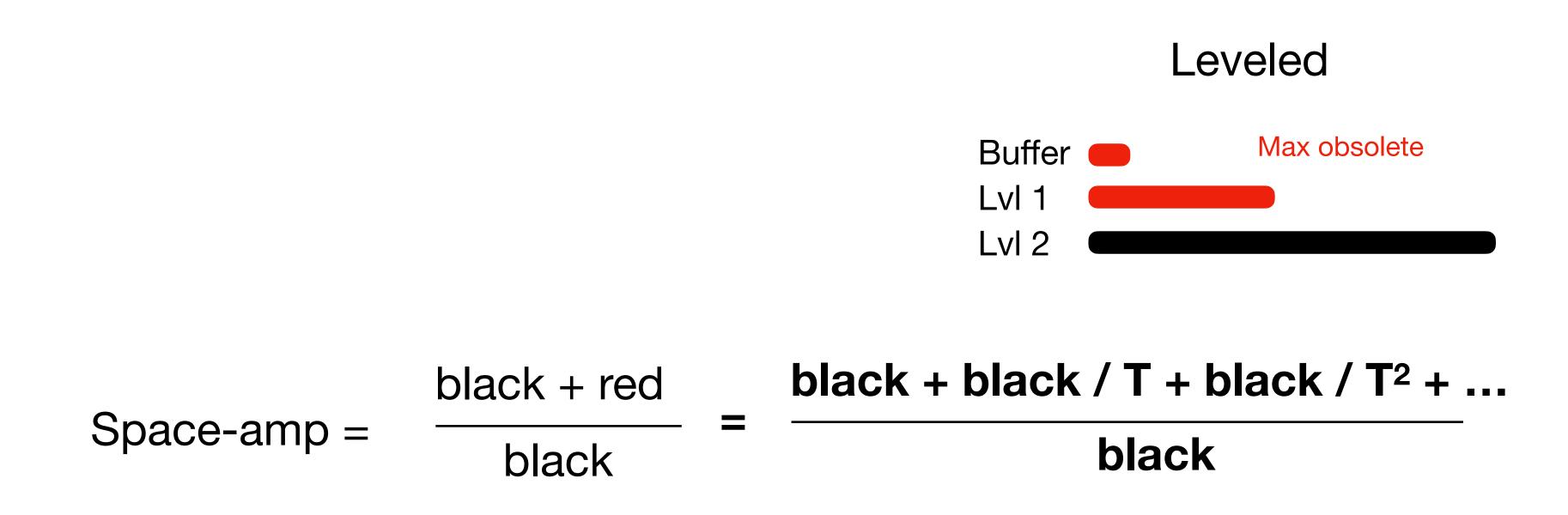
	Clustered	Unclustered	Unclustered is better when
#levels:	L _C =log _T (N/B)	Lu=log _T (N/K)	Large data entries (i.e., K > B)Few long range queries (S is small)
put:	O(L _C / B)	$O(1/B + L_U/K)$	- Gets are not a clear cut
get:	O(L _C)	O(L _U + 1)	
scan:	O(L _C +S/B)	O(L _U +S)	







$$Space-amp = \frac{black + red}{black}$$





Space-amp =
$$\frac{\text{black + red}}{\text{black}}$$
 = $\frac{\text{black + black / T + black / T^2 + ...}}{\text{black}}$ = $1 + \frac{1}{T} + \frac{1}{T^2} + ... \leq \frac{T}{T-1}$

