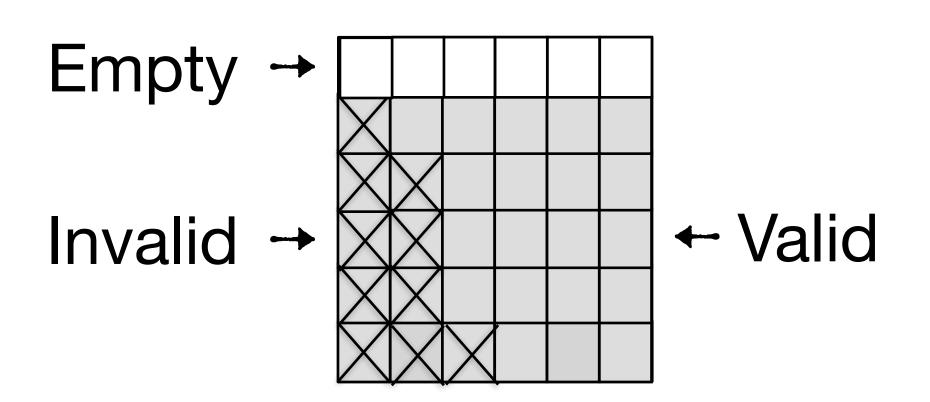
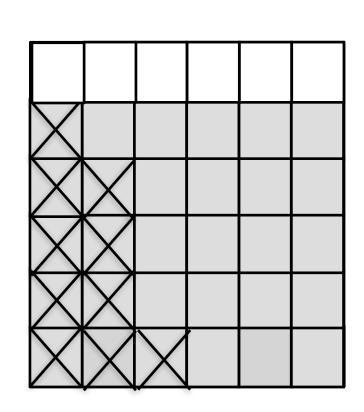
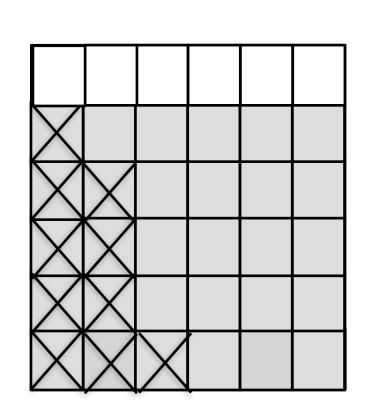
Indexing Tutorial

Database System Technology





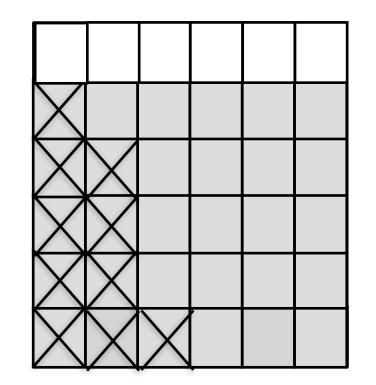


Consider an SSD with logical address space size L and physical capacity P. The usable capacity is a fraction of L/P of the overall capacity.

Worst case write-amplification (WA) occurs when each erase-unit has same number of invalid pages.

Worst case

Non-Worst Case

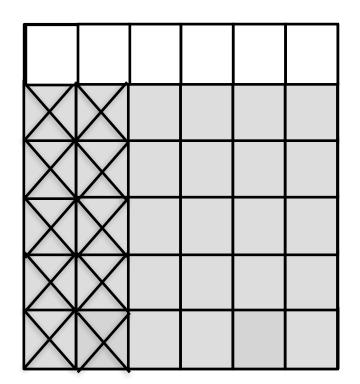


Best target

Consider an SSD with logical address space size L and physical capacity P. The usable capacity is a fraction of L/P of the overall capacity.

Worst case write-amplification (WA) occurs when each erase-unit has same number of invalid pages.

Worst case

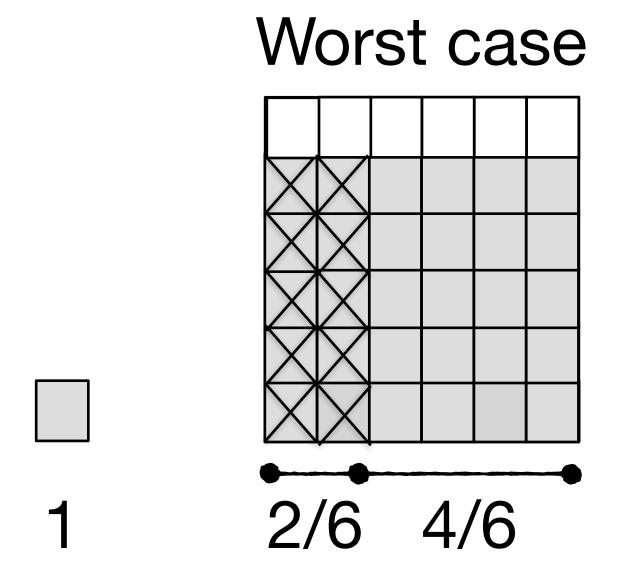


SSD WA Approx

$$1 + \frac{x}{1 - x}$$

Consider an SSD with logical address space size L and physical capacity P. The usable capacity is a fraction of L/P of the overall capacity.

Worst case write-amplification (WA) occurs when each erase-unit has same number of invalid pages.



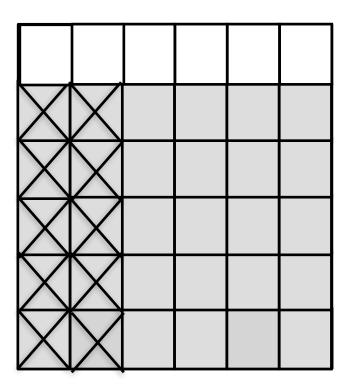
SSD WA Approx

$$1+\frac{x}{1-x}$$

Consider an SSD with logical address space size L and physical capacity P. The usable capacity is a fraction of L/P of the overall capacity.

Worst case write-amplification (WA) occurs when each erase-unit has same number of invalid pages.

Worst case



$$1 + \frac{4/6}{2/6}$$

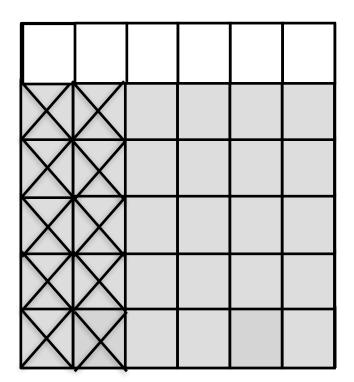
SSD WA Approx

$$1 + \frac{x}{1 - x}$$

Consider an SSD with logical address space size L and physical capacity P. The usable capacity is a fraction of L/P of the overall capacity.

Worst case write-amplification (WA) occurs when each erase-unit has same number of invalid pages.

Worst case



 $1 + \frac{4}{2} = 3$

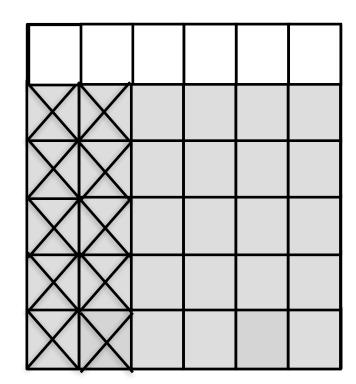
SSD WA Approx

$$1 + \frac{x}{1 - x}$$

Consider an SSD with logical address space size L and physical capacity P. The usable capacity is a fraction of L/P of the overall capacity.

Worst case write-amplification (WA) occurs when each erase-unit has same number of invalid pages.

Worst case



SSD WA Approx

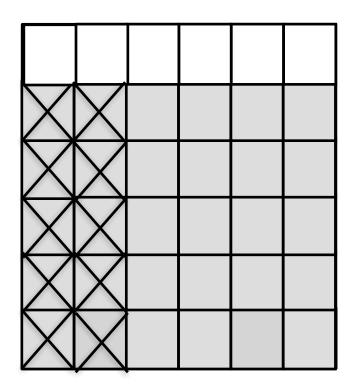
$$1 + \frac{x}{1 - x}$$

In worst case: x = valid / total = 4/6 = L/P

Consider an SSD with logical address space size L and physical capacity P. The usable capacity is a fraction of L/P of the overall capacity.

Worst case write-amplification (WA) occurs when each erase-unit has same number of invalid pages.

Worst case

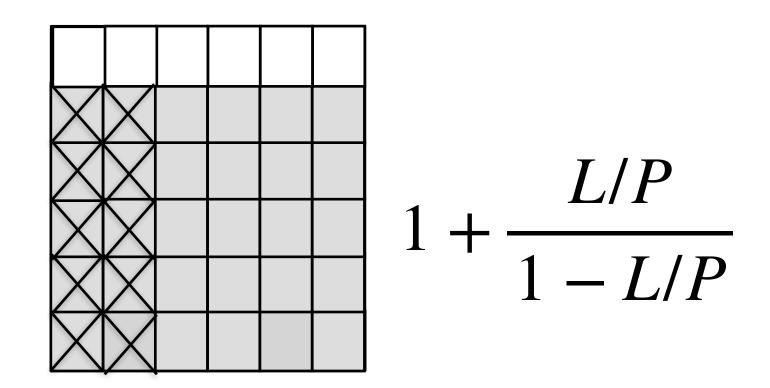


SSD WA Approx

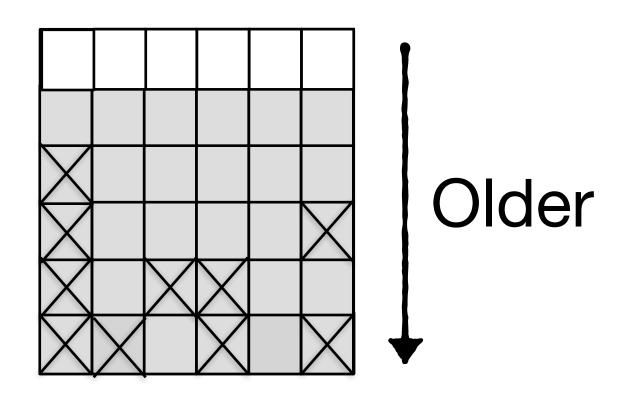
$$1 + \frac{L/P}{1 - L/P}$$

In worst case: x = valid / total = 4/6 = L/P

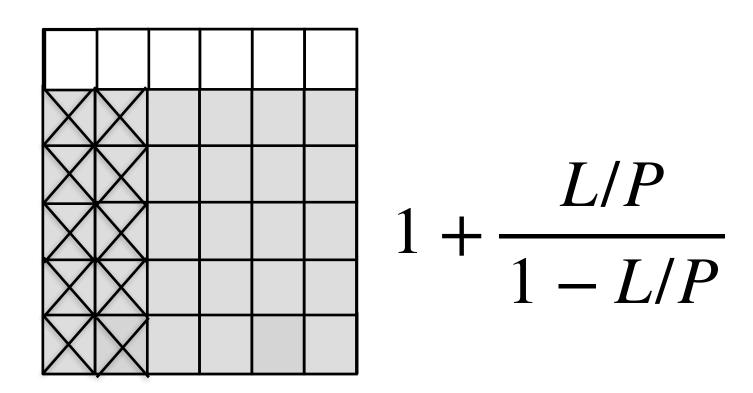
But the worst-case hardly ever happens in practice

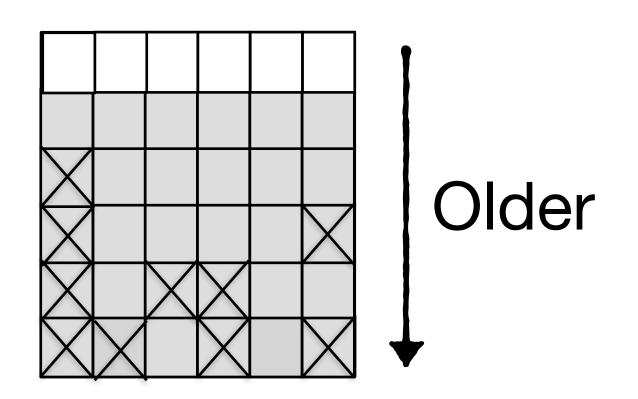


In practice, erase units written longer ago have more invalid pages, so GC is cheaper

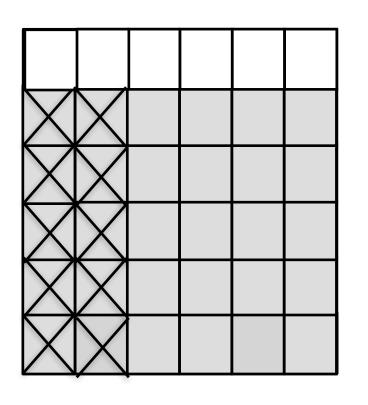


Cost model assuming uniformly randomly distributed writes?



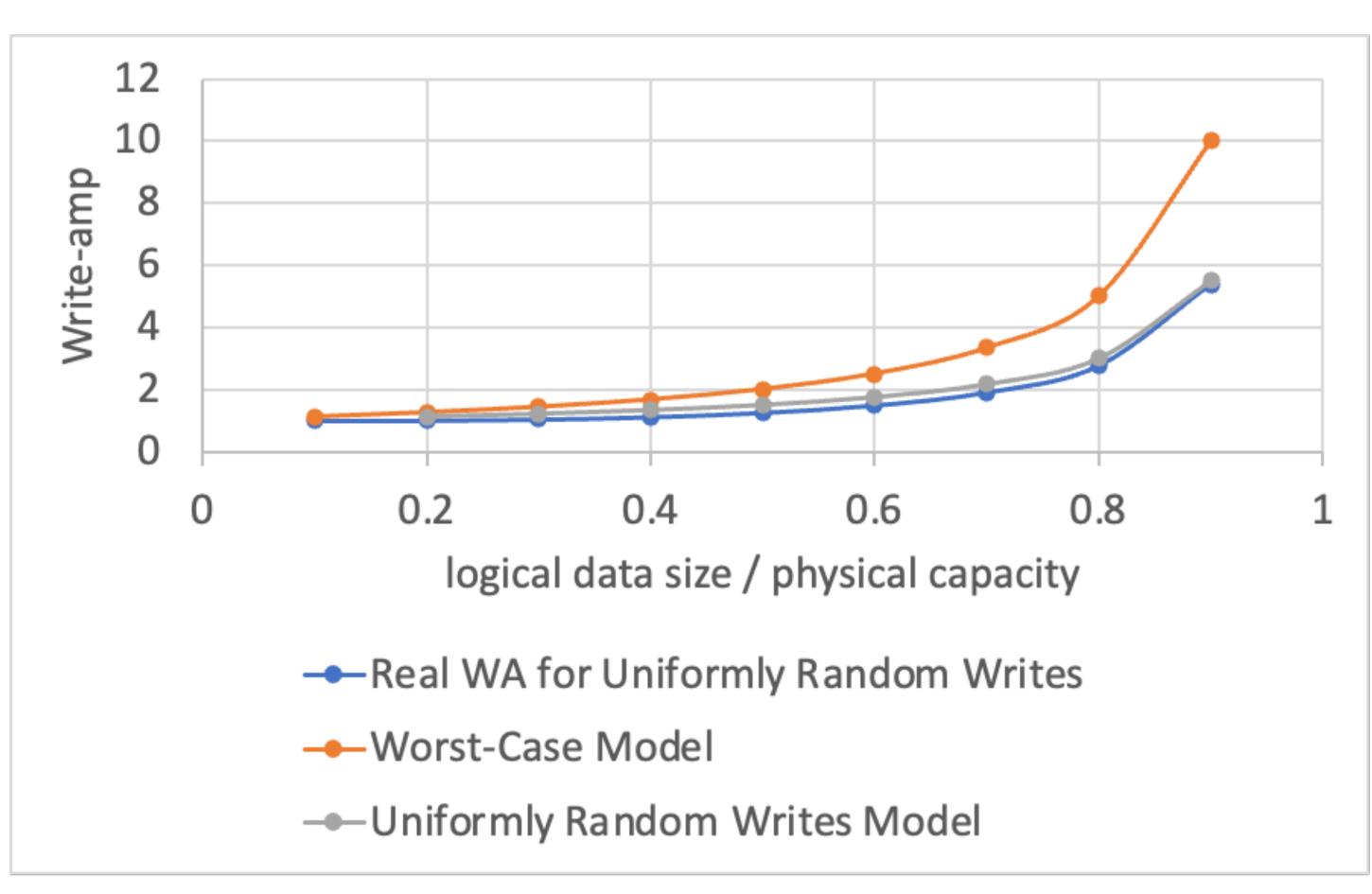


Cost model assuming uniformly randomly distributed writes?



$$1 + \frac{L/P}{1 - L/P}$$

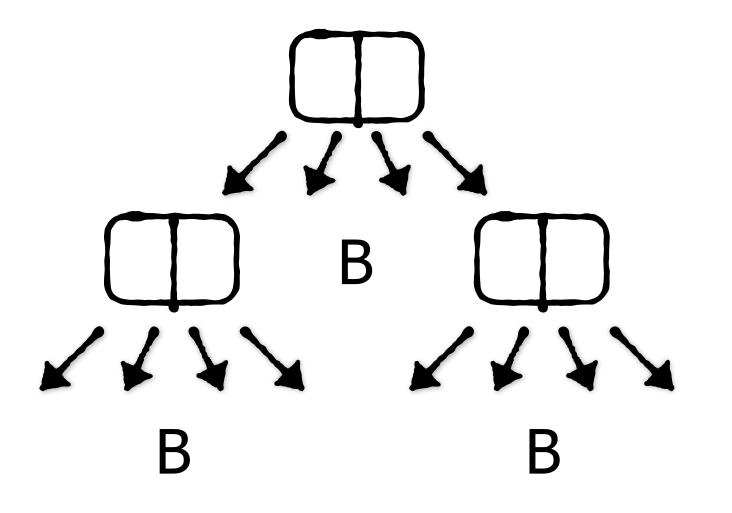
$$1 + \frac{1}{2} \cdot \frac{L/P}{1 - L/P}$$



$$1 + \frac{L/P}{1 - L/P}$$

$$1 + \frac{1}{2} \cdot \frac{L/P}{1 - L/P}$$

Consider a B-tree subject to uniformly randomly distributed updates. There are 100 entries per page. The b-tree occupies 70% of the physical SSD capacity, while the rest is over-provisioned. What write-amplification would you expect? A back-of-the-envelope calculation is enough.





Consider a B-tree subject to uniformly randomly distributed updates. There are 100 entries per page. The b-tree occupies 70% of the physical SSD capacity, while the rest is over-provisioned. What write-amplification would you expect? A back-of-the-envelope calculation is enough.

SSD WA Model with uniformly random updates: $1 + \frac{1}{2} \cdot \frac{L/P}{1 - I/P}$

B-tree update cost: ≈ 1 write I/O

Consider a B-tree subject to uniformly randomly distributed updates. There are 100 entries per page. The b-tree occupies 70% of the physical SSD capacity, while the rest is over-provisioned. What write-amplification would you expect? A back-of-the-envelope calculation is enough.

SSD WA Model with uniformly random updates: $1 + \frac{1}{2} \cdot \frac{L/P}{1 - I/P}$

$$1 + \frac{1}{2} \cdot \frac{L/P}{1 - L/P}$$

B-tree update cost: ≈ 1 write I/O

B-tree WA: ≈ B

Consider a B-tree subject to uniformly randomly distributed updates. There are 100 entries per page. The b-tree occupies 70% of the physical SSD capacity, while the rest is over-provisioned. What write-amplification would you expect? A back-of-the-envelope calculation is enough.

SSD WA Model with uniformly random updates: $1 + \frac{1}{2} \cdot \frac{L/P}{1 - I/P}$

B-tree update cost: ≈ 1 write I/O B-tree WA: ≈ B

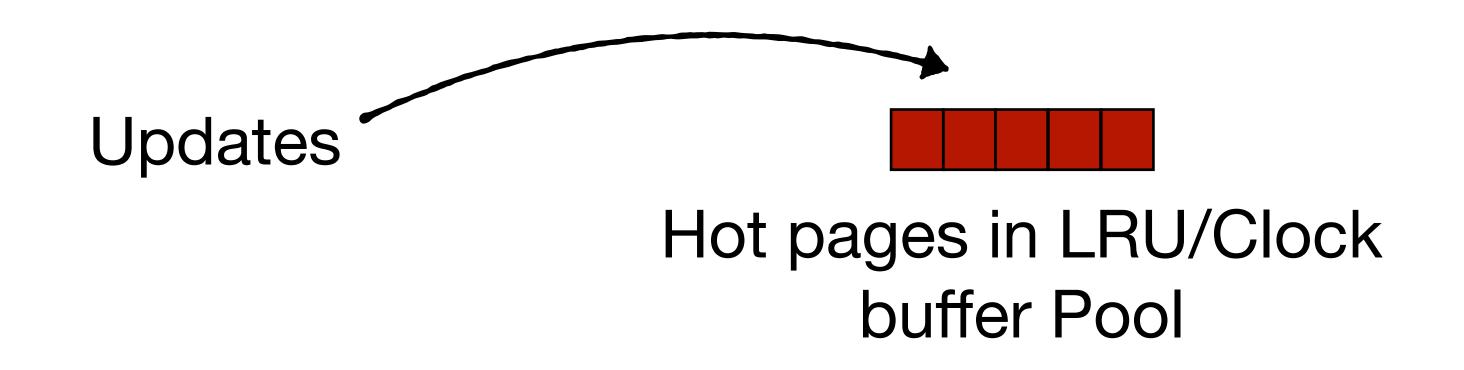
Total write-amp: $100 \cdot (1 + \frac{1}{2} \cdot \frac{0.7}{1 - 0.7}) = 216.67$

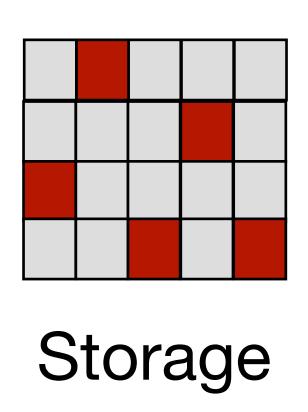
Consider a B-tree subject to uniformly randomly distributed updates. There are 100 entries per page. The b-tree occupies 70% of the physical SSD capacity, while the rest is over-provisioned. What write-amplification would you expect? A back-of-the-envelope calculation is enough.

Now suppose the workload exhibits skew (some entries are more likely to be updated). Which mechanism of a database allows us to reduce write-amplification in this case, and why?

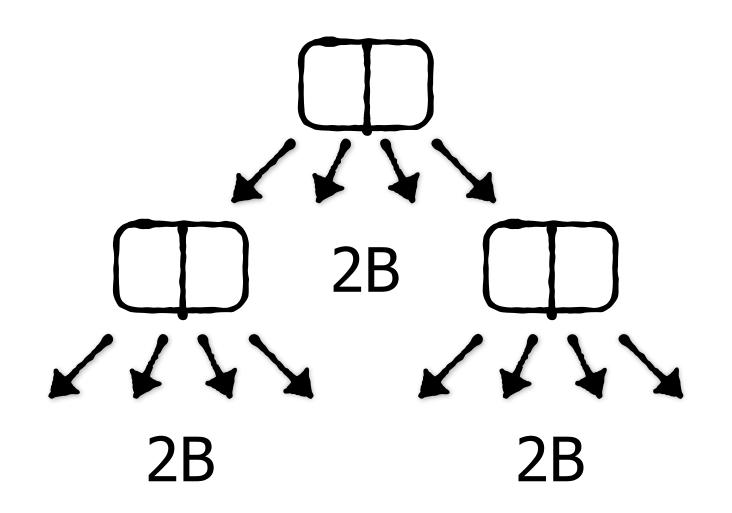
Consider a B-tree subject to uniformly randomly distributed updates. There are 100 entries per page. The b-tree occupies 70% of the physical SSD capacity, while the rest is over-provisioned. What write-amplification would you expect? A back-of-the-envelope calculation is enough.

Now suppose the workload exhibits skew (some entries are more likely to be updated). Which mechanism of a database allows us to reduce write-amplification in this case, and why?





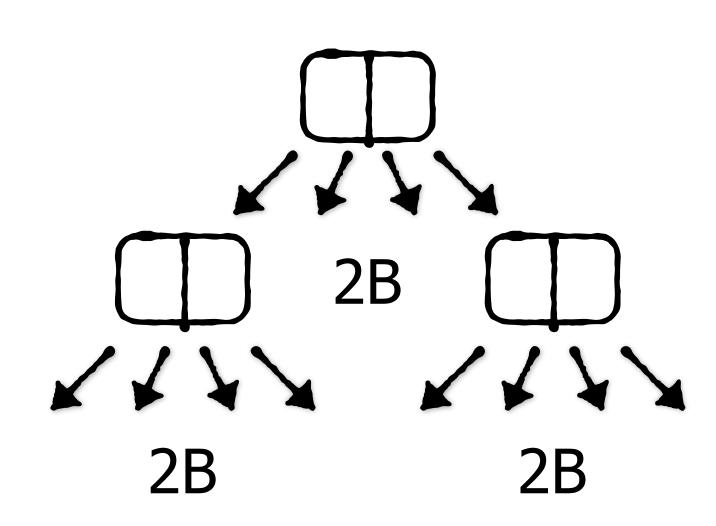
We have a B-tree on an SSD. the workload is read-intensive, so we only care about read performance. Each node and each SSD page are 4 KB. Consider making each B-tree node take up two rather than one flash pages (i.e., 8KB rather than 4KB). This can make the tree shallower. Is this a good idea?



We have a B-tree on an SSD. the workload is read-intensive, so we only care about read performance. Each node and each SSD page are 4 KB. Consider making each B-tree node take up two rather than one flash pages (i.e., 8KB rather than 4KB). This can make the tree shallower. Is this a good idea?

On SSD, cost is measured as # pages accessed

$$= 2 \cdot \log_{2B}(N)$$



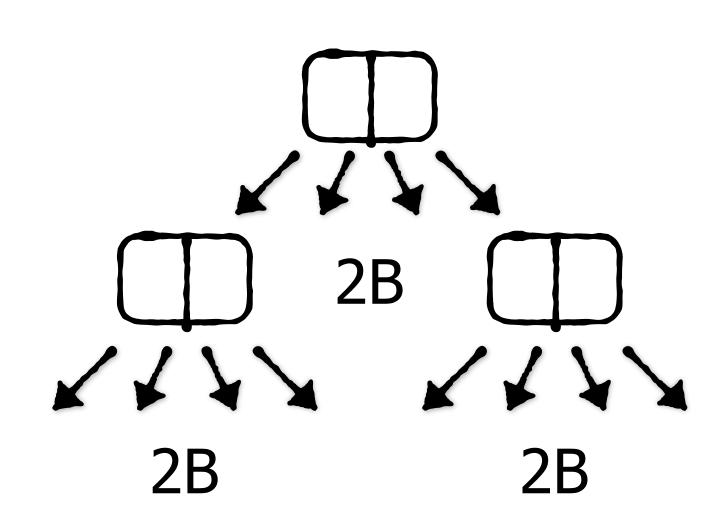
We have a B-tree on an SSD. the workload is read-intensive, so we only care about read performance. Each node and each SSD page are 4 KB. Consider making each B-tree node take up two rather than one flash pages (i.e., 8KB rather than 4KB). This can make the tree shallower. Is this a good idea?

On SSD, cost is measured as # pages accessed

 $2 \cdot \log_{2B}(N) \leq \log_{B}(N)$

1

Condition for being cheaper than standard B-tree

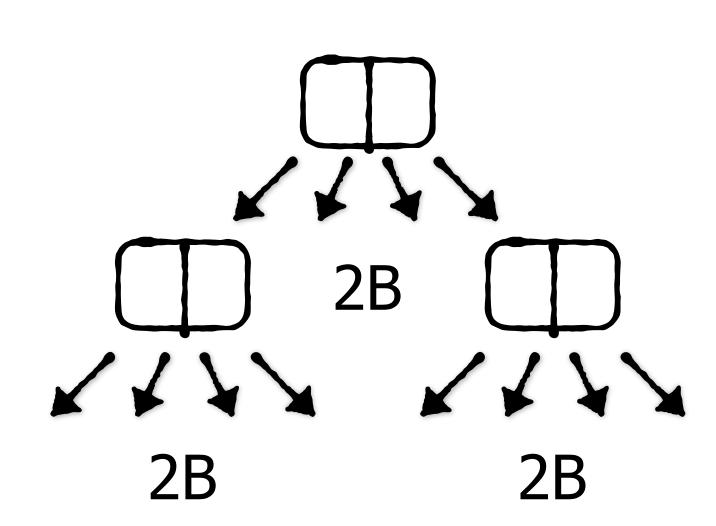


We have a B-tree on an SSD. the workload is read-intensive, so we only care about read performance. Each node and each SSD page are 4 KB. Consider making each B-tree node take up two rather than one flash pages (i.e., 8KB rather than 4KB). This can make the tree shallower. Is this a good idea?

On SSD, cost is measured as # pages accessed

$$2 \cdot \log_{2B}(N) \leq \log_{B}(N)$$

Simplifies to: B ≤ 2



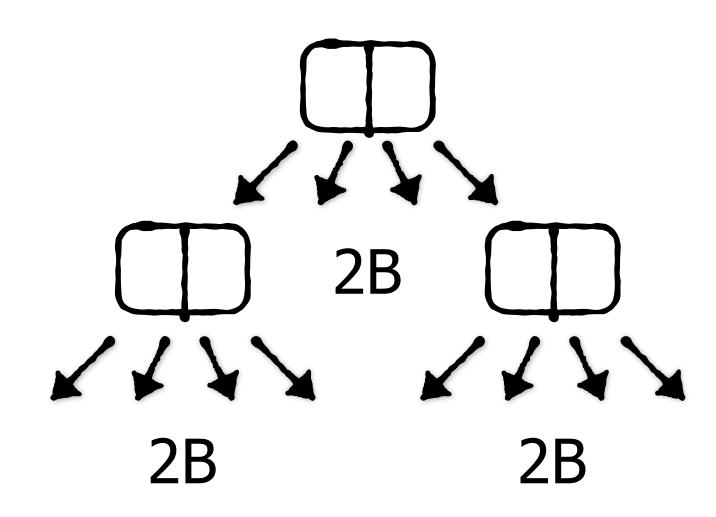
We have a B-tree on an SSD. the workload is read-intensive, so we only care about read performance. Each node and each SSD page are 4 KB. Consider making each B-tree node take up two rather than one flash pages (i.e., 8KB rather than 4KB). This can make the tree shallower. Is this a good idea?

On SSD, cost is measured as # pages accessed

$$2 \cdot \log_{2B}(N) \leq \log_{B}(N)$$

Simplifies to: $B \le 2$

But B is typically larger. So it's not generally a good idea.



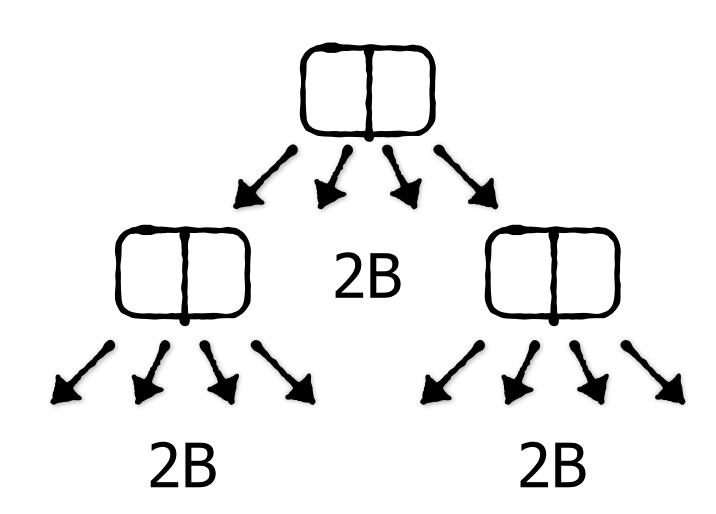
We have a B-tree on an SSD. the workload is read-intensive, so we only care about read performance. Each node and each SSD page are 4 KB. Consider making each B-tree node take up two rather than one flash pages (i.e., 8KB rather than 4KB). This can make the tree shallower. Is this a good idea?

On SSD, cost is measured as # pages accessed

$$2 \cdot \log_{2B}(N) \leq \log_{B}(N)$$

Simplifies to: $B \le 2$

What if we use disk instead of SSD?



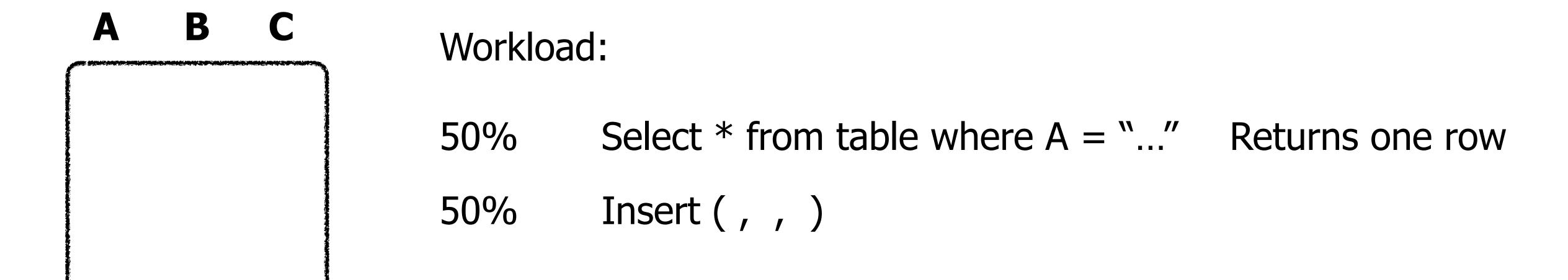
We have a B-tree on an SSD. the workload is read-intensive, so we only care about read performance. Each node and each SSD page are 4 KB. Consider making each B-tree node take up two rather than one flash pages (i.e., 8KB rather than 4KB). This can make the tree shallower. Is this a good idea?

What if we use disk instead of SSD?

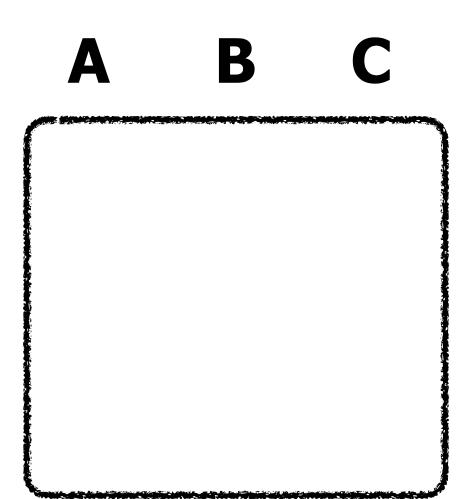


On disk, seek & rotational delay dominates, while data transfer is negligible. So this is a good idea. A node size in the area of a few MB is common for disk. Beyond that, we begin to incur a penalty for long sequential accesses.

Consider a table with columns A, B and C. Suppose we employ buffered insertions at a cost of O(1/B) each.



Consider a table with columns A, B and C. Suppose we employ buffered insertions at a cost of O(1/B) each.



Workload:

50% Select * from table where A = "..." Returns one row

50% Insert (,,)

Should we employ a B-tree index on any of the columns? Estimate the overall I/O cost of both operations with & without it. (Back-of-the-envelope reasoning is sufficient)

Consider a table with columns A, B and C. Suppose we employ buffered insertions at a cost of O(1/B) each.

A B C

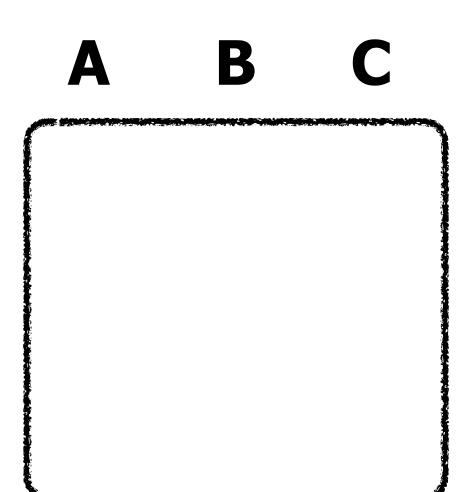
Select * from table where A = "..." Returns one row

50% Insert (,,)

Should we employ a B-tree index on any of the columns? Estimate the overall I/O cost of both operations with & without it.

Query Insert
$$I/O$$
 cost without index $0.5 \cdot N/B + 0.5 \cdot 1/B \approx 0.5 \cdot N/B$

Consider a table with columns A, B and C. Suppose we employ buffered insertions at a cost of O(1/B) each.



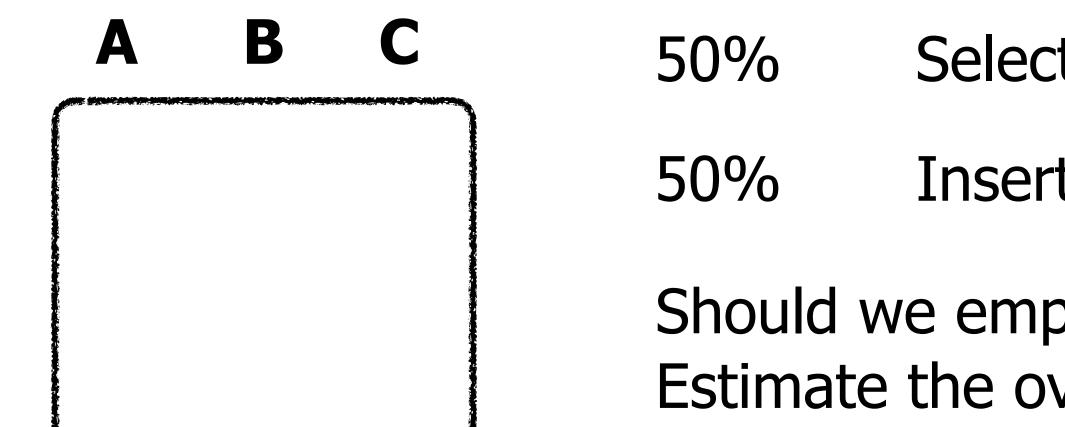
50% Select * from table where A = "..." Returns one row

50% Insert (,,)

Should we employ a B-tree index on any of the columns? Estimate the overall I/O cost of both operations with & without it.

	Query	Insert	
I/O cost without index	0.5 · N/B +	0.5 · 1/B	$\approx 0.5 \cdot N/B$
I/O cost with index on A	0.5 · log _B N +	0.5 · (1/B +	$-\log_B N) \approx \log_B N$

Consider a table with columns A, B and C. Suppose we employ buffered insertions at a cost of O(1/B) each.



Select * from table where A = "..." Returns one row

Insert (,,)

Should we employ a B-tree index on any of the columns? Estimate the overall I/O cost of both operations with & without it.

Insert Query I/O cost without index 0.5 · N/B I/O cost with index on A $0.5 \cdot \log_B N + 0.5 \cdot (1/B + \log_B N) \approx \log_B N$

Indexing column A significantly reduces overall costs.

A table has columns A, B and C. The workload consists of two query types:

A B C

50% Select A from table where A = "..."

Returns one row

50% Select * from table where B > x and B < y Returns avg. 10 rows

How should we index this table? B-tree or hash table? Clustered vs. unclustered? Estimate worst-case I/O cost with your plan for each query with these indexes assuming $N=2^{40}$ and $B=2^{10}$

A table has columns A, B and C. The workload consists of two query types:

A B C

50% Select A from table where A = "..."

Returns one row

50% Select * from table where B > x and B < y

Returns avg. 10 rows

How should we index this table? B-tree or hash table? Clustered vs. unclustered? Estimate worst-case I/O cost with your plan for each query with these indexes assuming $N=2^{40}$ and $B=2^{10}$

Clustered B-tree on B

Unclustered Hash table on A

A table has columns A, B and C. The workload consists of two query types:

A B C

50% Select A from table where A = "..."

Returns one row

50% Select * from table where B > x and B < y

Returns avg. 10 rows

How should we index this table? B-tree or hash table? Clustered vs. unclustered? Estimate worst-case I/O cost with your plan for each query with these indexes assuming $N=2^{40}$ and $B=2^{10}$

Clustered B-tree on B

Unclustered Hash table on A

 $\lceil \log_{B}(N) \rceil + \lceil S/B \rceil$ 4 + 1

≈1 I/O per query