## LSM-Trees

Niv Dayan



## **Project Announcements**



**Feedback weeks** 

Oct 6-10 - Step 1

Nov 3-7 - Step 2



**Deadline** 

Dec 1

## **Project Announcements**



Min Expected Contribution: 20%



Struggling in your group? Ask for help.

## **Project Announcements**



Min Expected Contribution: 20%

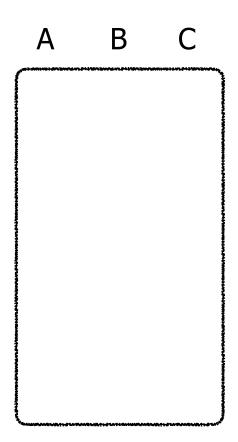


Struggling in your group? Ask for help.



Last resort: splitting

## We are trying to optimize the following kinds of queries



select \* from table where A = "..."

select \* from table where B > "..." and B < "..."

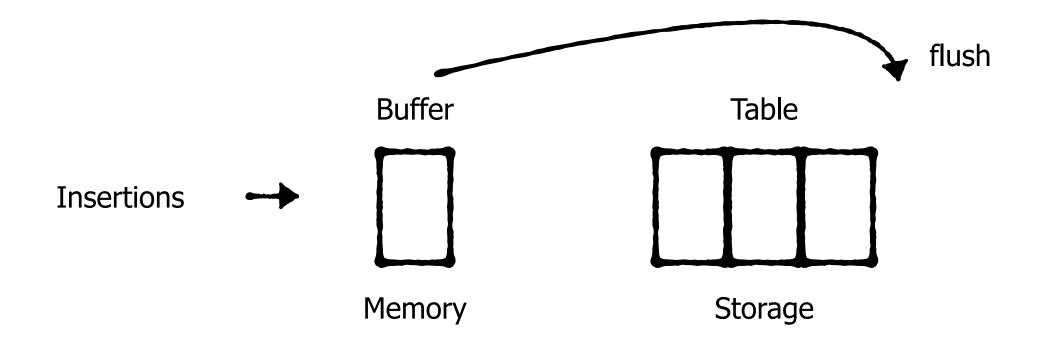
## Recap

Append-Only Table Sorted Table Extendible Hashing B-tree

## **Append-Only Table**

O(N/B) for any selection query

O(1/B) for insertions



Append-Only Table Sorted Table Extendible Hashing B-tree

#### Sorted Table

A ...

7 ...

22 ...

32 ...

32 ...

**61** ...

**-**

74 ...

90 ...

**97** ...

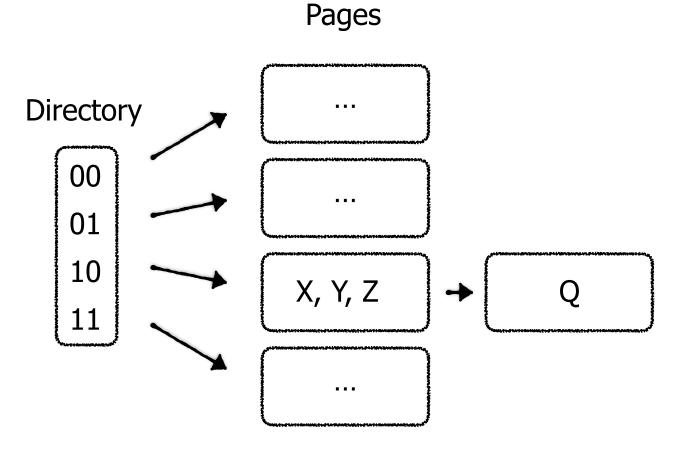
Binary search: O(log<sub>2</sub> N/B) I/O

Update/insert/delete: O(N/B) read & write I/O

Append-Only Table Sorted Table Extendible Hashing B-tree

A directory maps pages in storage with a given hash prefix

Handle overflows via chaining

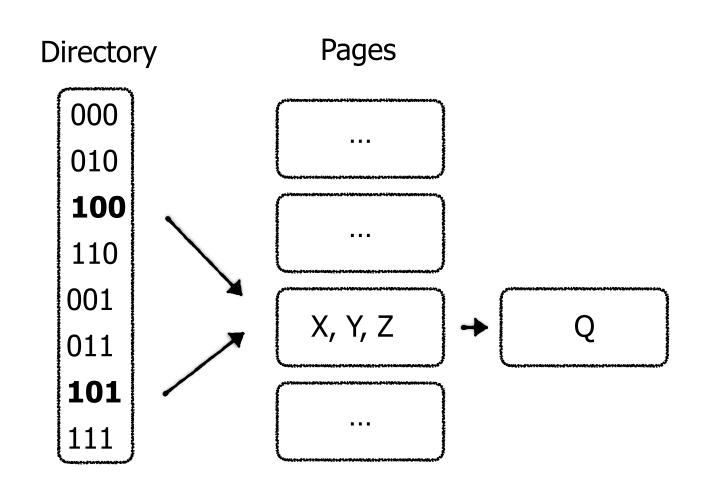


Memory

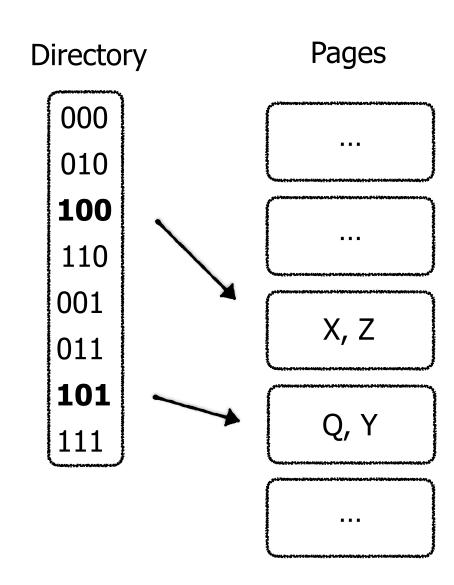
Storage

When we reach capacity, double directory size.

New directory slots still point to previous pages



Split one overflowing bucket at a time by rehashing entries.

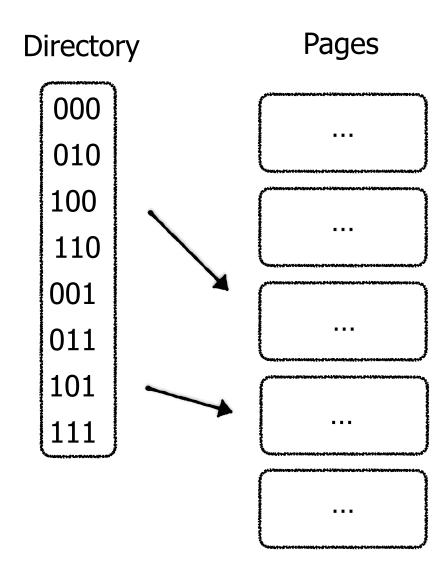


Query cost: O(1) read I/O

Insertion cost:

O(1) read I/O

O(1) write I/O



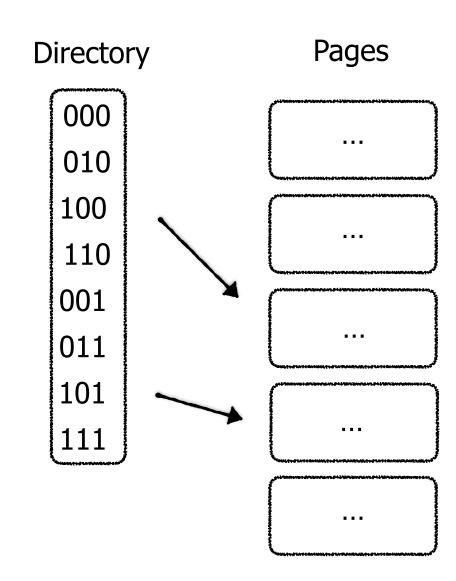
Query cost: O(1) read I/O

Insertion cost:

O(1) read I/O

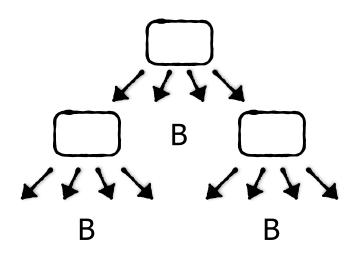
O(1) \* **GC** write I/O

Random writes to storage entail SSD garbage-collection



Append-Only Table Sorted Table Extendible Hashing B-Tree

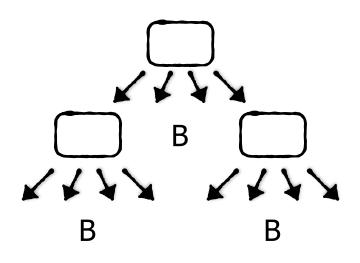
#### **B-Tree**



Each node has B children. This allows pruning by a factor of B in each level

Can issue random writes to storage leading to SSD garbage-collection

#### **B-Tree**



Each node has B children. This allows pruning by a factor of B in each level

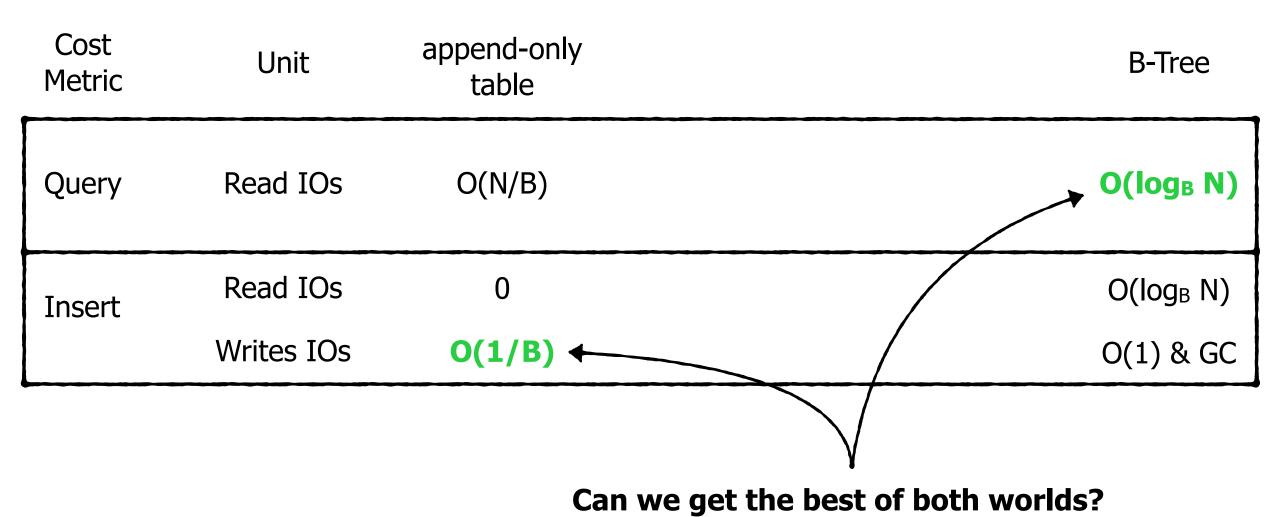
Can issue random writes to storage leading to SSD garbage-collection

**Query:** O(log<sub>B</sub> N) read I/O

**Update/insert/query:** O(log<sub>B</sub> N) read I/O

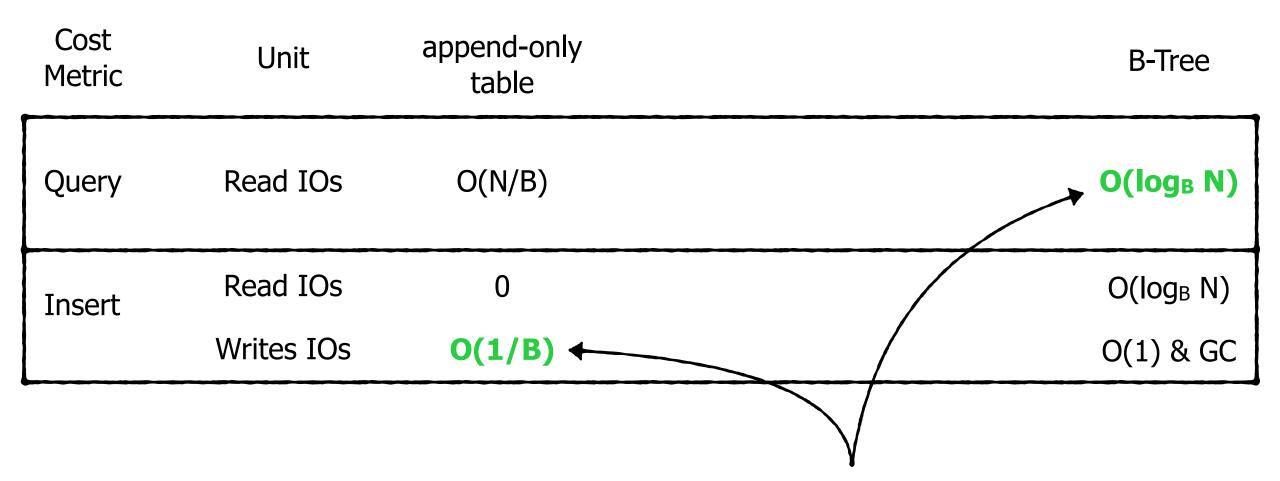
& O(1) \* GC write I/O

Cost Metric	Unit	append-only table	Sorted File	Extendible Hashing	B-Tree
Query	Read IOs	O(N/B)	O(log <sub>2</sub> N/B)	O(1)	O(log <sub>B</sub> N)
Insert	Read IOs	0	O(N/B)	O(1)	O(log <sub>B</sub> N)
	Writes IOs	O(1/B)	O(N/B)	O(1) & GC	O(1) & GC



Typical CS approach: (1) Identify solutions with diff properties

(2) Combine to get something in-between



Can we get the best of both worlds?

## The Log-Structured Merge-Tree





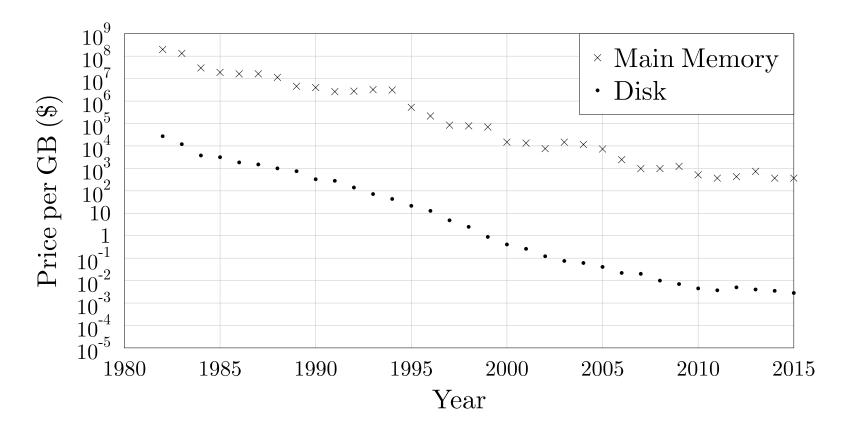
B-tree	LSM-tree is	Google's BigTable	LSM-tree is
is invented	invented	adopts LSM-tree	widely used
1970	1996	2005	Today



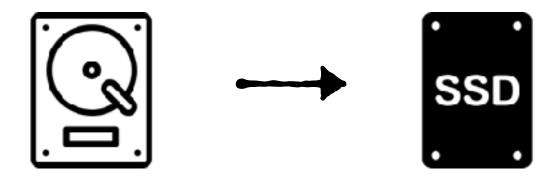
B-tree is invented	LSM-tree is invented	Google's BigTable adopts LSM-tree	LSM-tree is widely used
1970	1996	2005	Today

Why was it not invented and used sooner?

The declining costs of storage allow us to store more data cheaply. Hence, application workloads are becoming more write-intensive. This drives a need to optimize for data ingestion.



At the same, the advent of SSDs makes write I/Os more expensive than read I/Os. The drives a need for more write-optimized data structures.



# Necessity is the mother of invention



Plato (Loose attribution)

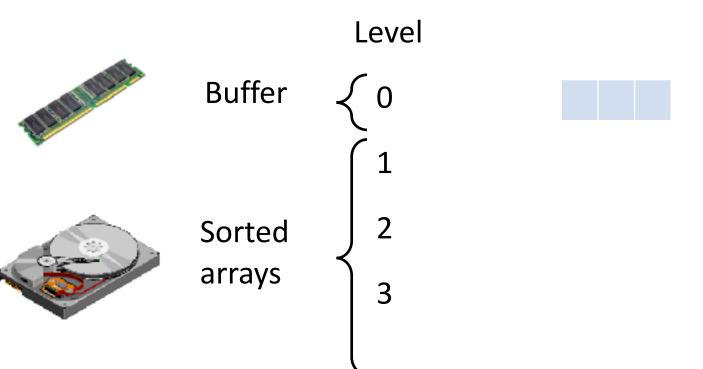
Leveled LSM-tree

Basic LSM-tree

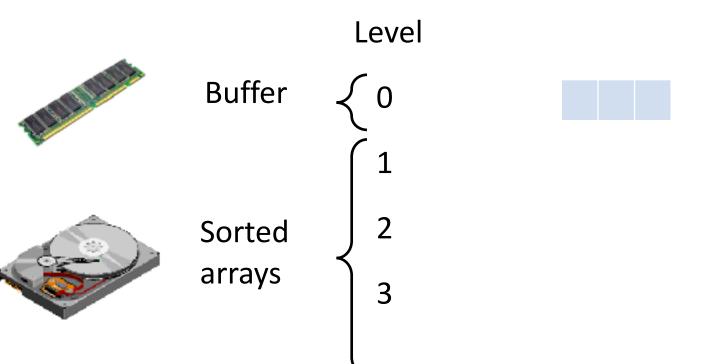
Tiered LSM-tree



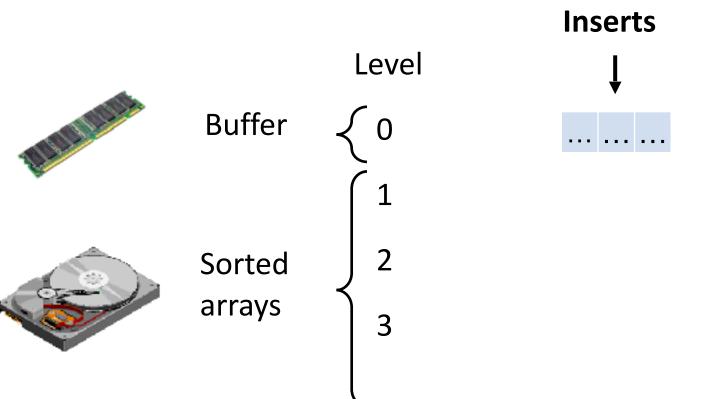




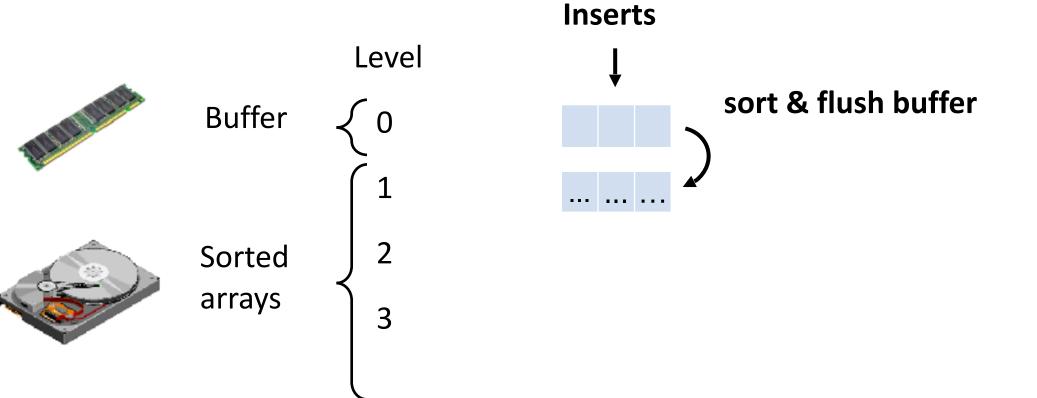
Design principle #1:



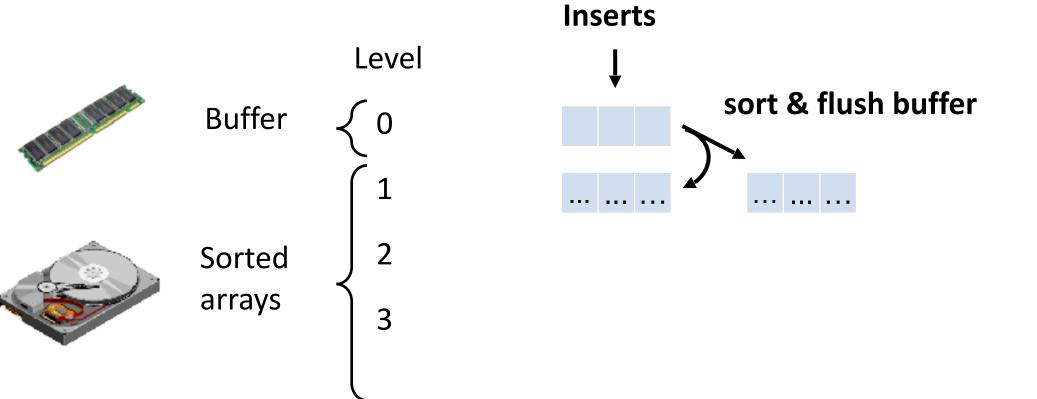
Design principle #1:



Design principle #1:

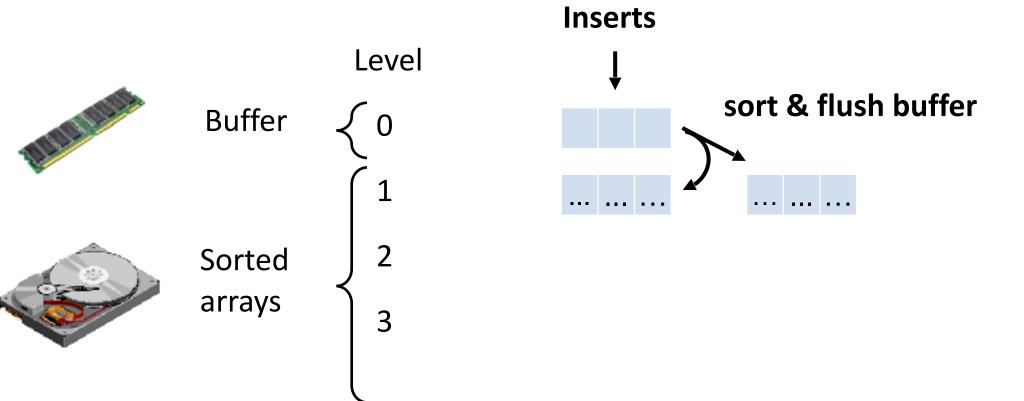


Design principle #1:



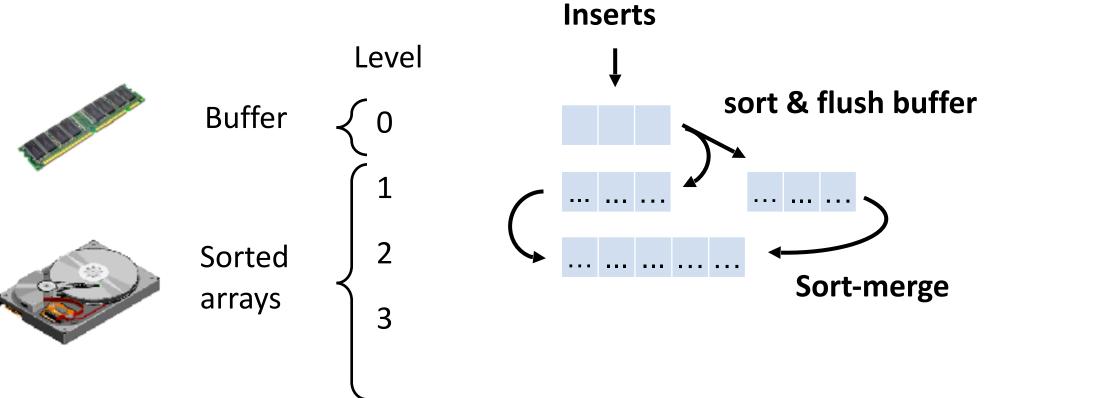
Design principle #1: optimize for insertions by buffering

Design principle #2: optimize for lookups by sort-merging SSTs



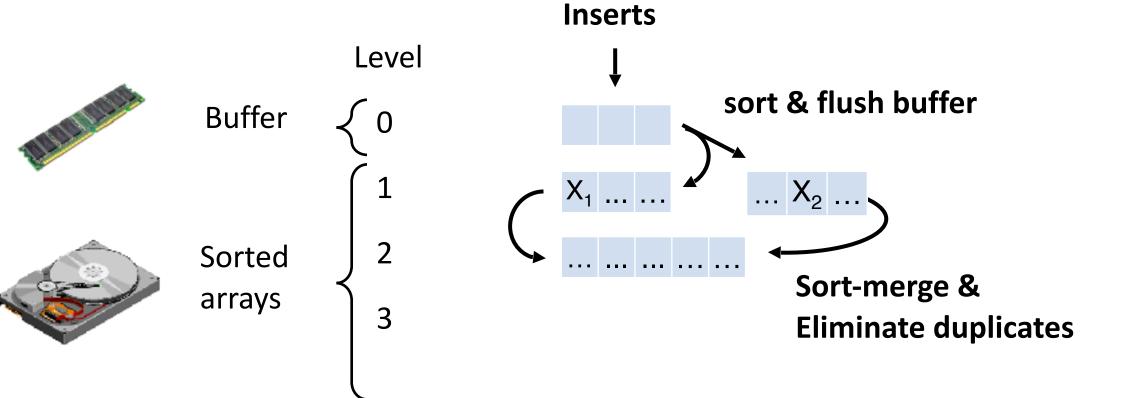
Design principle #1: optimize for insertions by buffering

Design principle #2: optimize for lookups by sort-merging SSTs



Design principle #1: optimize for insertions by buffering

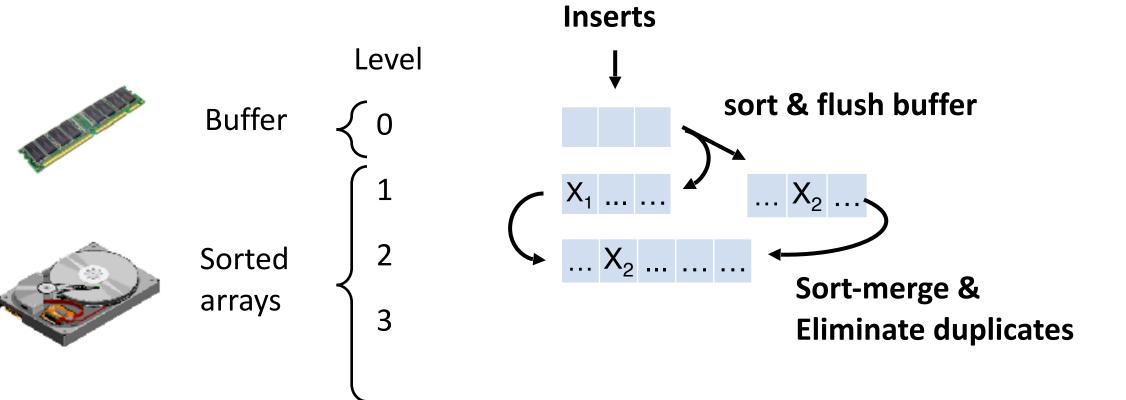
Design principle #2: optimize for lookups by sort-merging SSTs



#### **Basic LSM-tree**

Design principle #1: optimize for insertions by buffering

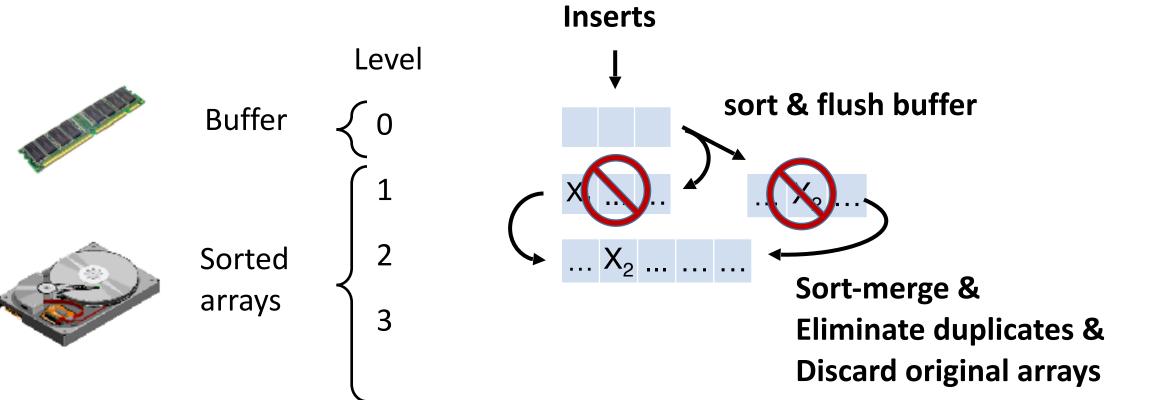
Design principle #2: optimize for lookups by sort-merging SSTs

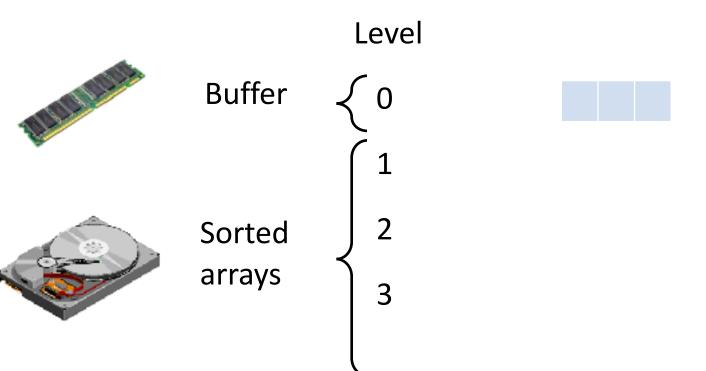


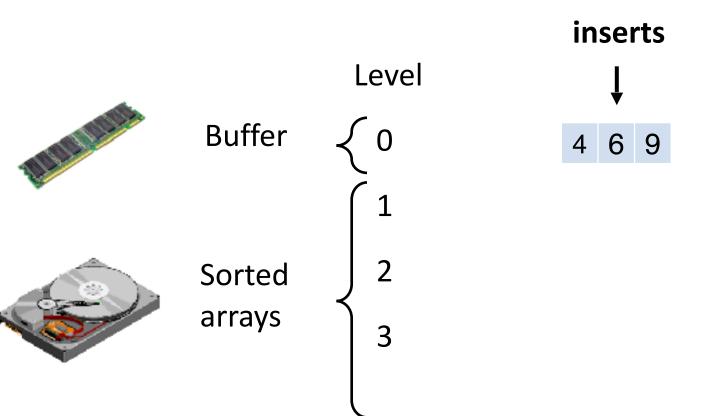
#### **Basic LSM-tree**

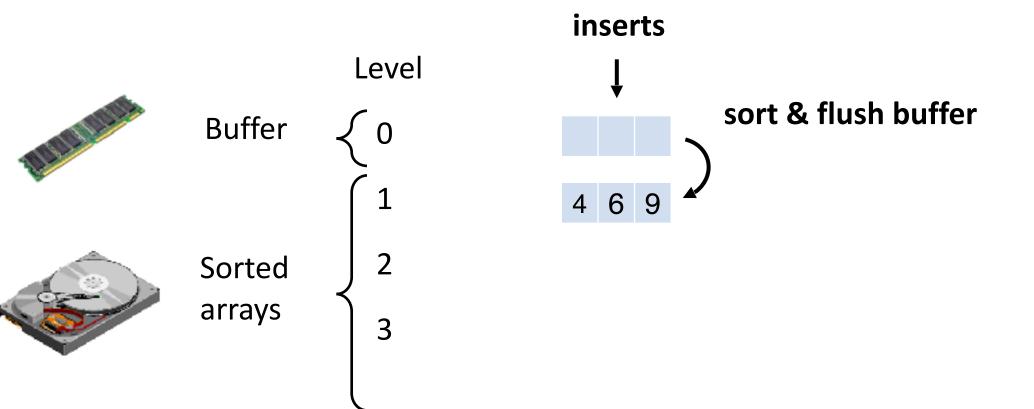
Design principle #1: optimize for insertions by buffering

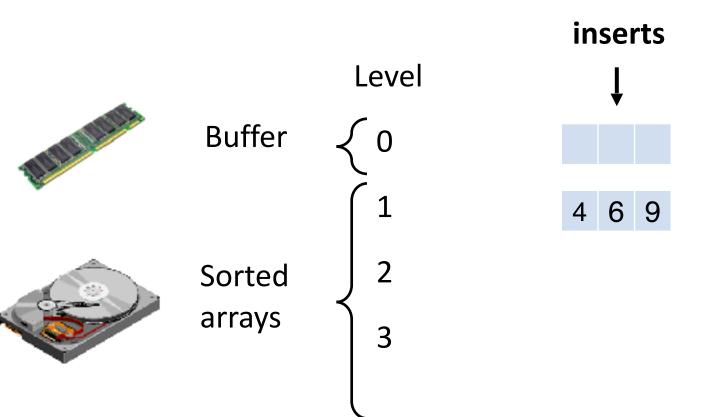
Design principle #2: optimize for lookups by sort-merging SSTs

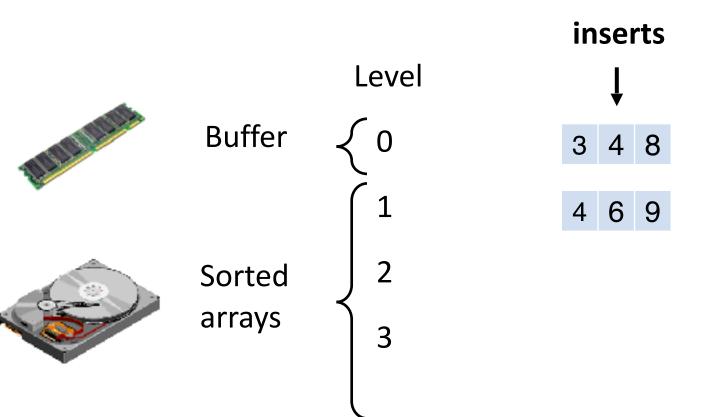


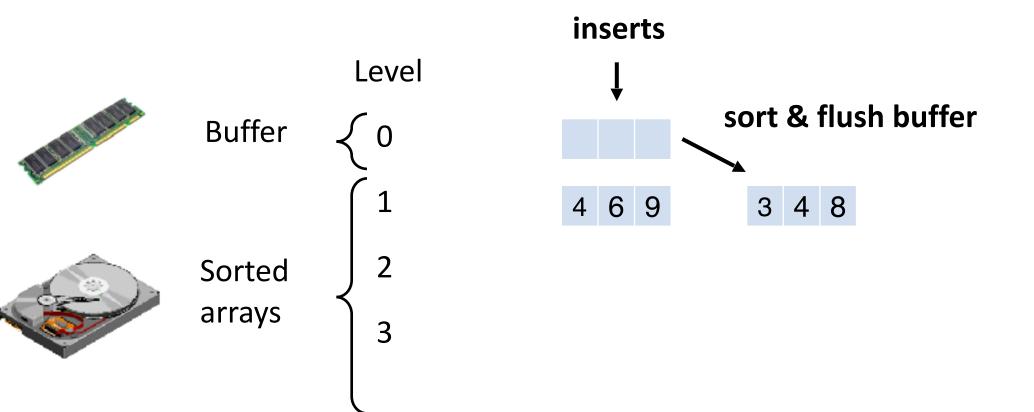


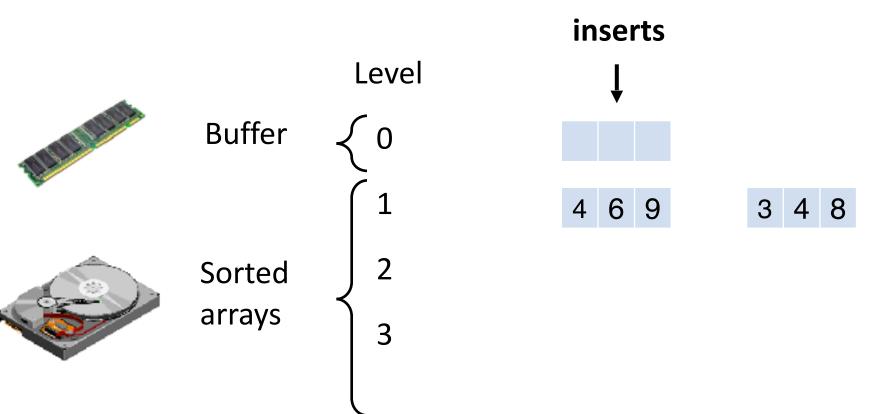


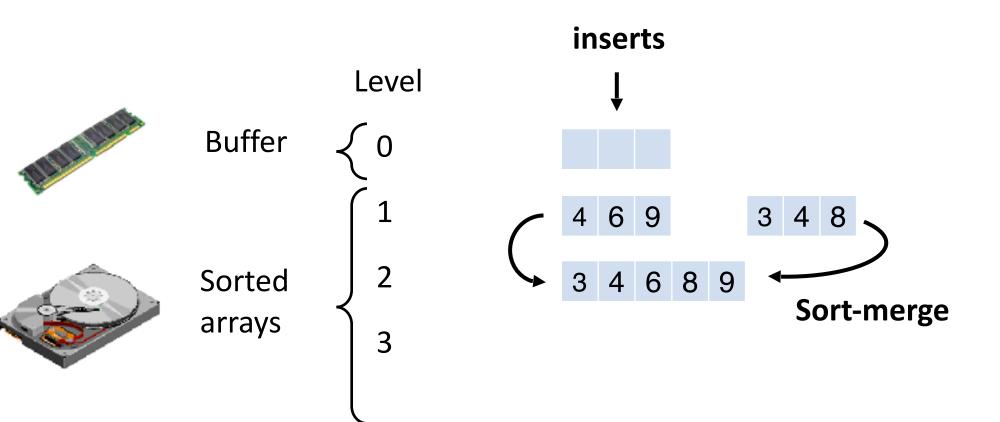


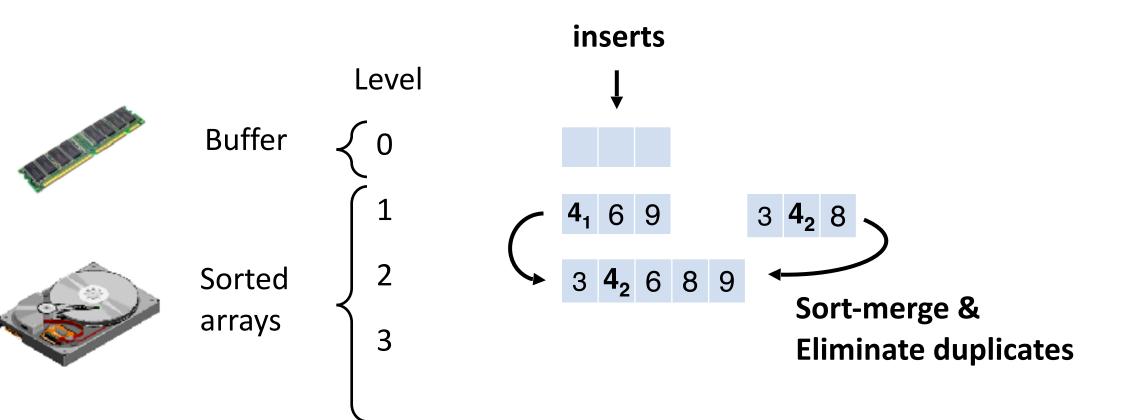


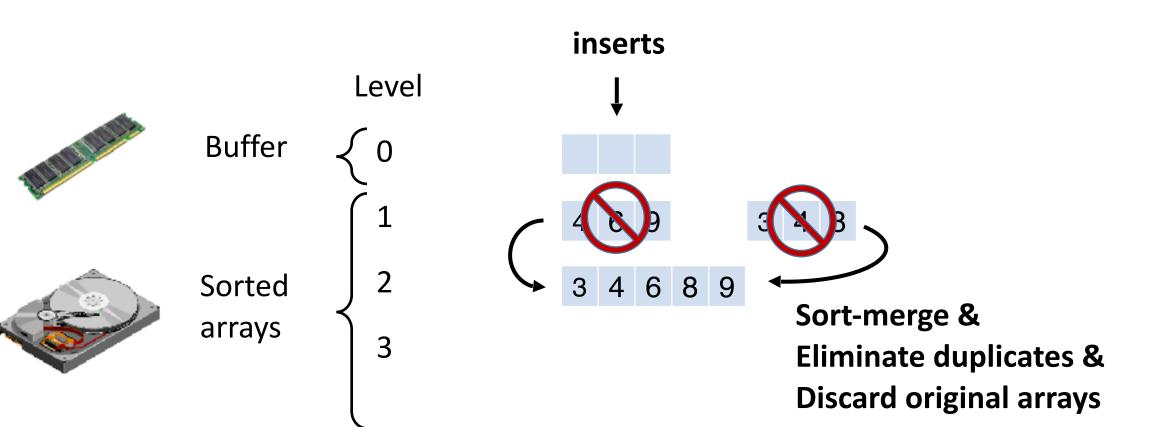


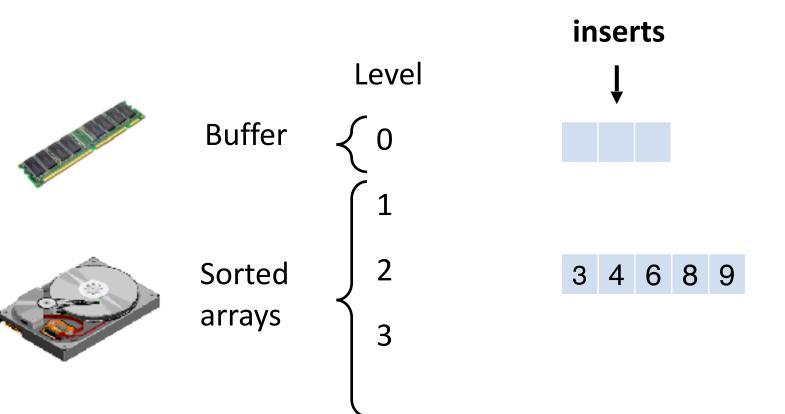


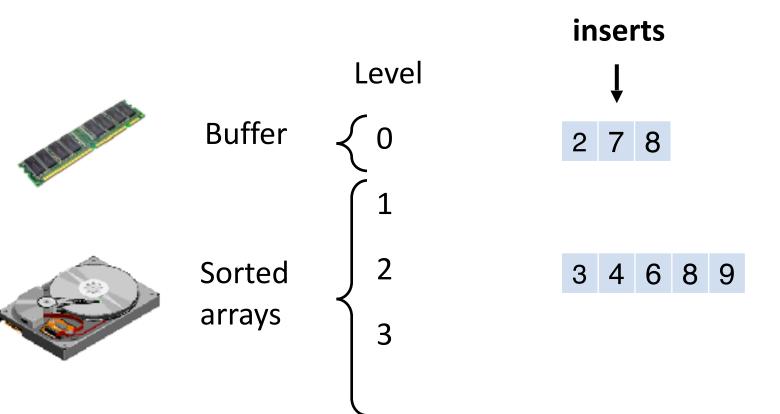


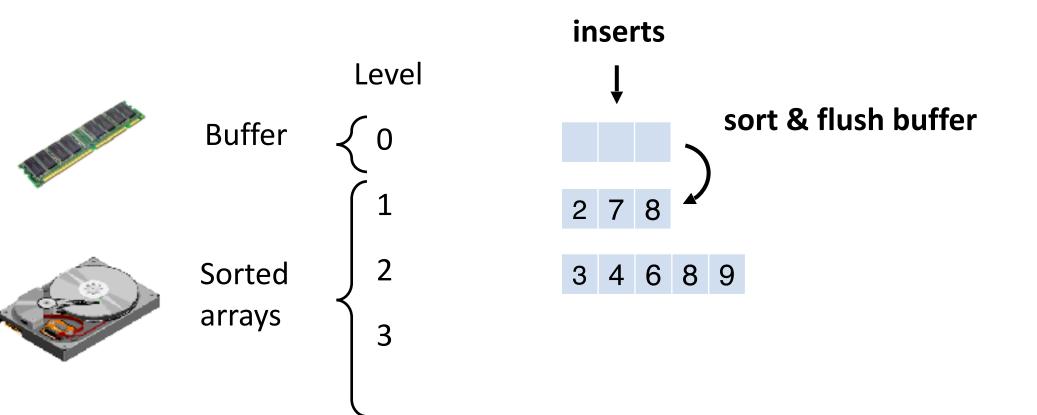


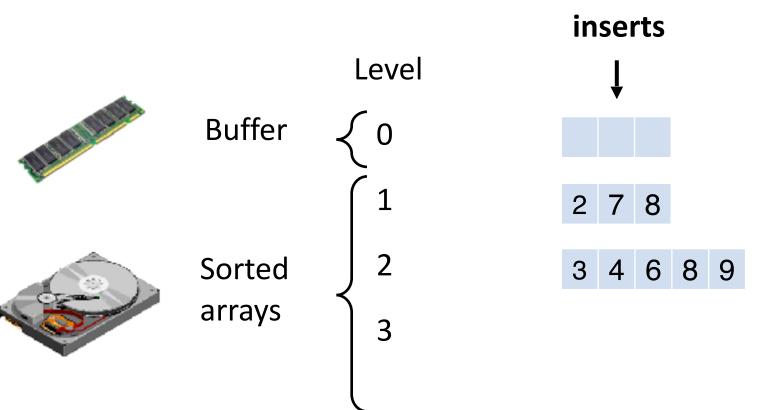




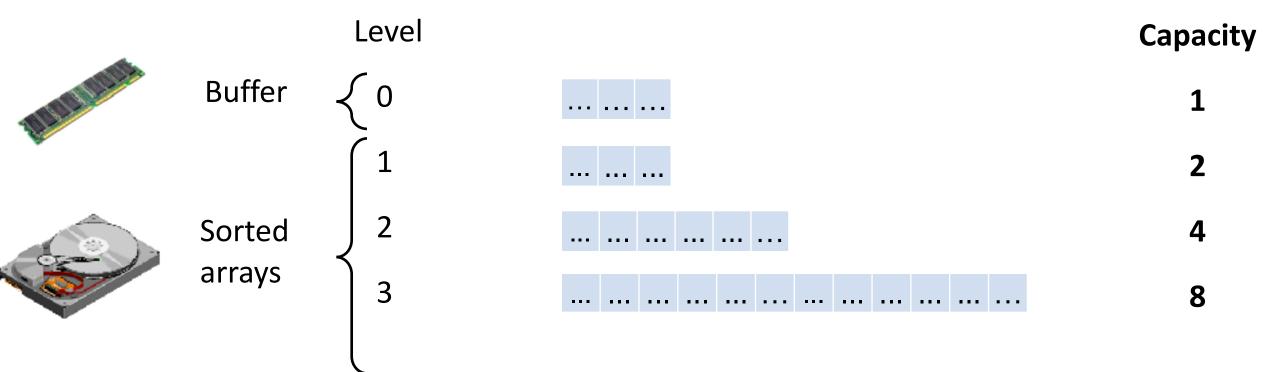




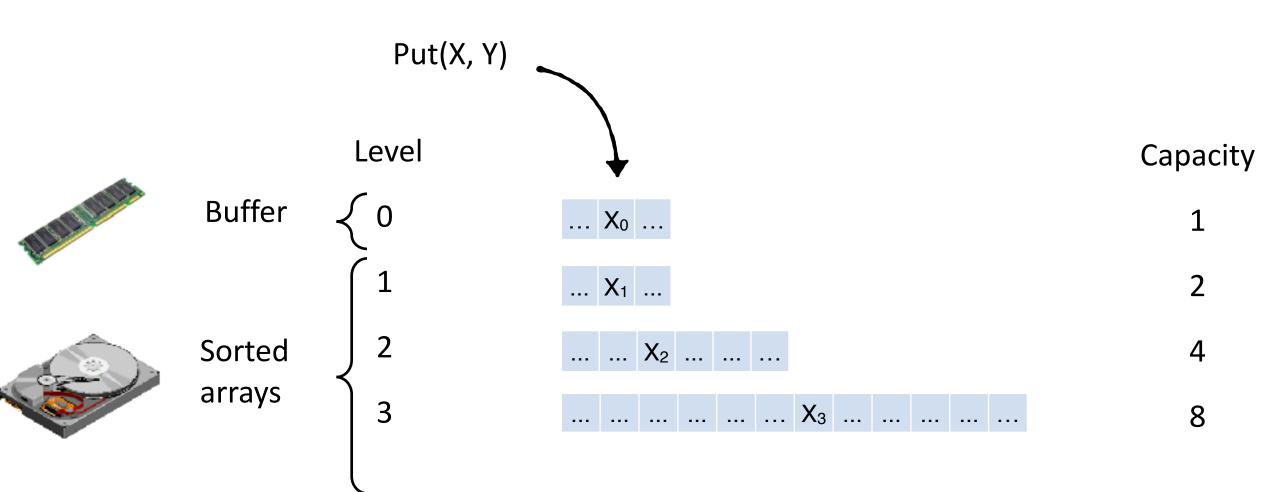




### Levels have exponentially increasing capacities.

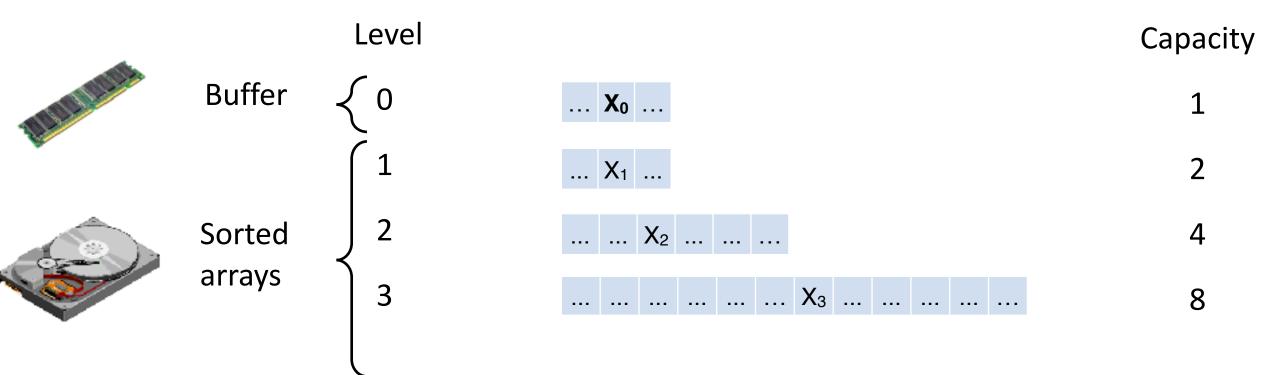


#### An updates is made out-of-place through an insertion into the buffer

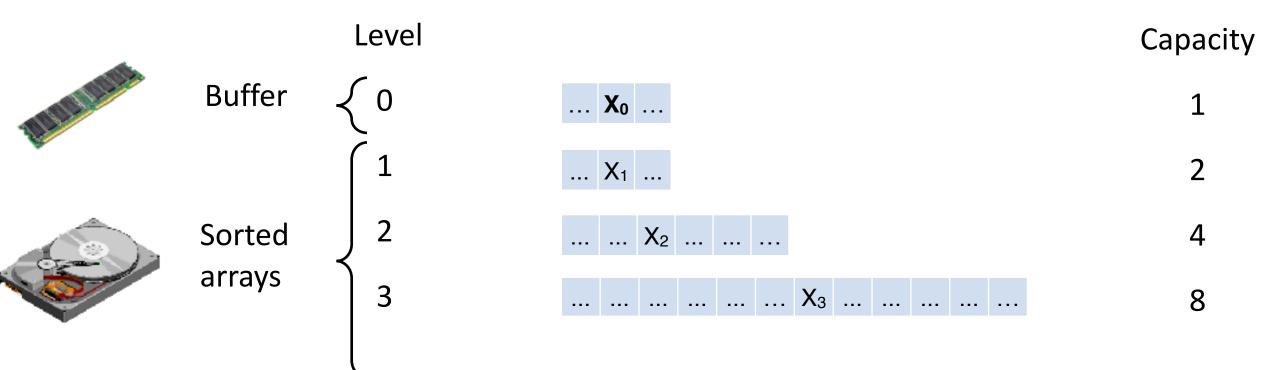


An updates is made out-of-place through an insertion into the buffer

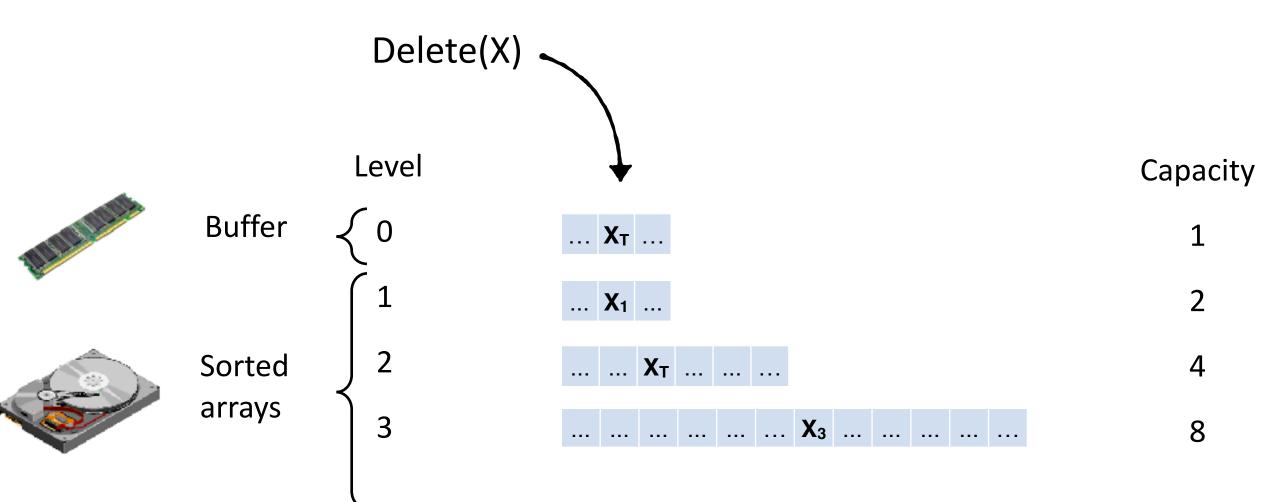
# Only an entry's most recent version should be returned to the user during a query



#### How to handle deletes?

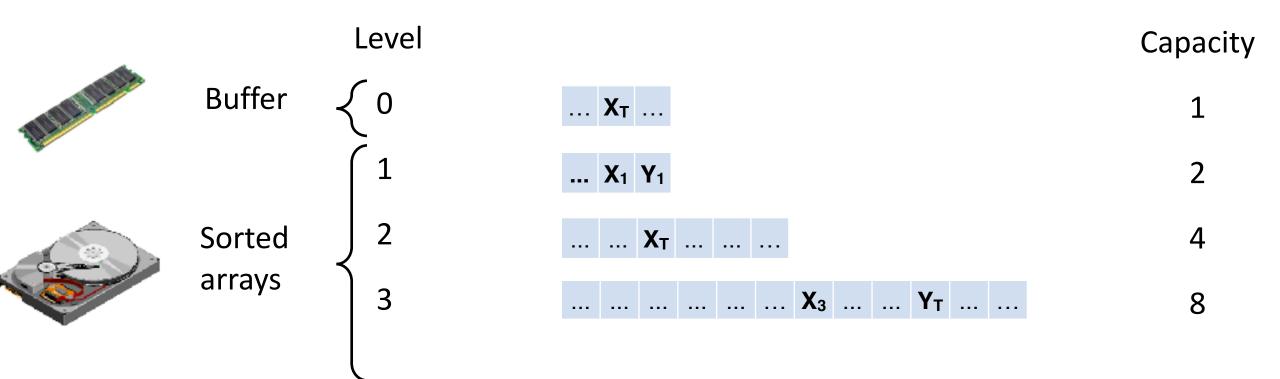


#### Delete an entry out-of-place by inserting a tombstone

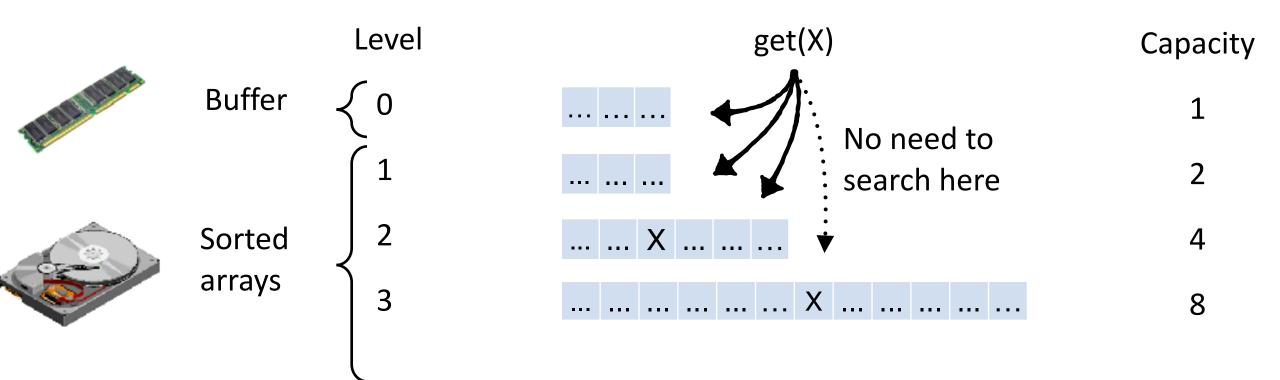


#### Delete an entry out-of-place by inserting a tombstone

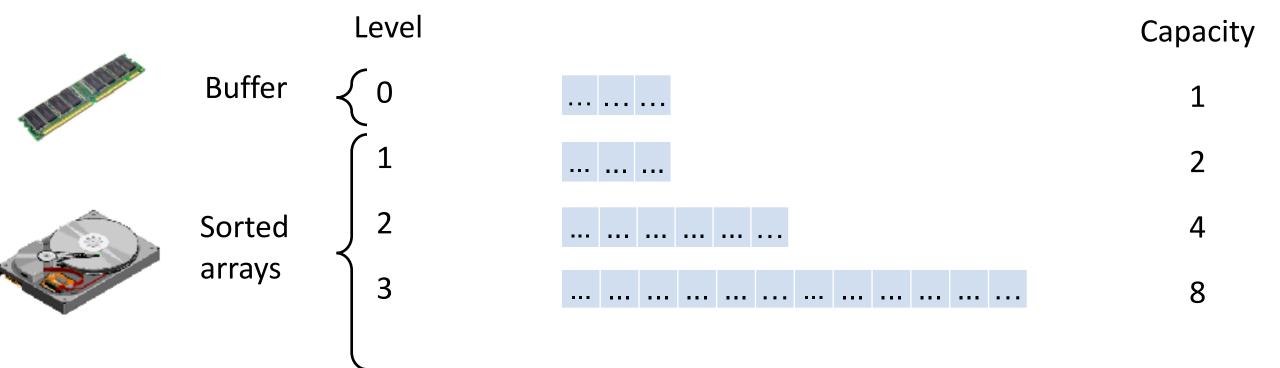
If the most recent version of an entry is a tombstone, it's considered deleted. (X is deleted but Y isn't in this example)



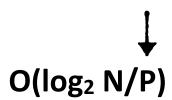
A point query searches the LSM-tree from smaller to larger levels. It stops when it finds the first entry with a matching key. Entries with a matching key at larger levels are older and thus superseded.



#### How many levels to search?



#### **Buffer size (# entries)**



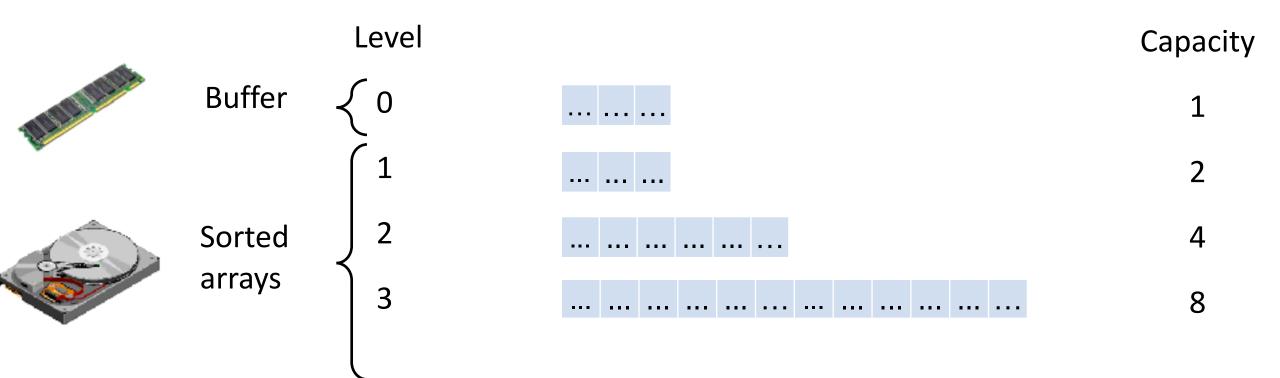
How many levels to search?

		Level	Capacity
William .	Buffer	€ 0	 1
•			 2
	Sorted arrays	2	 4
		3	 8

How many levels to search?

 $O(log_2 N/P)$ 

#### Cost per searching each level?

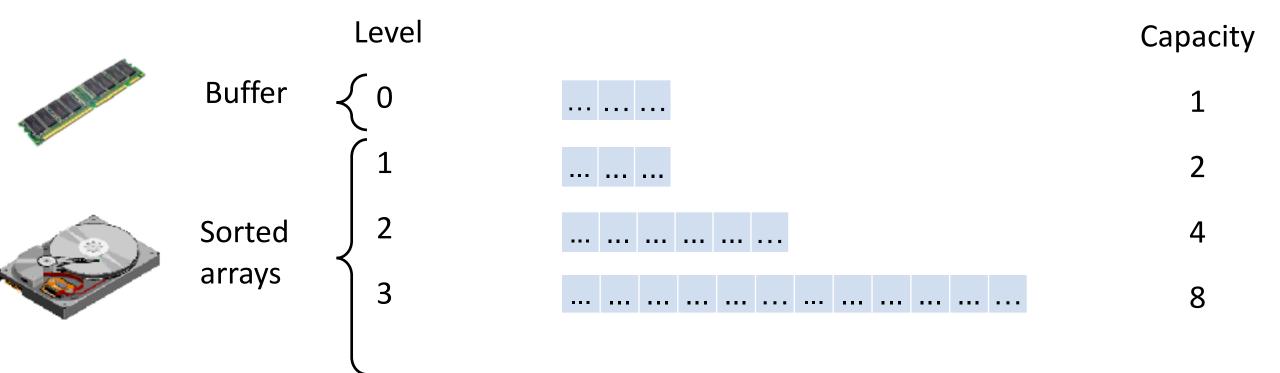


How many levels to search?

 $O(log_2 N/P)$ 

Cost per searching each level?

O(log<sub>2</sub> N/B) w. binary search



How many levels to search?  $O(log_2 N/P)$ 

Cost per searching each level? O(log<sub>2</sub> N/B)

Total search cost:  $O(log_2(N/P) * log_2(N/B))$ 

		Level		Capacity
	Buffer	€ 0		1
			*** *** ***	2
	Sorted arrays	2		4
		3		8

How many levels to search?

 $O(log_2 N/P)$ 

Cost per searching each level?

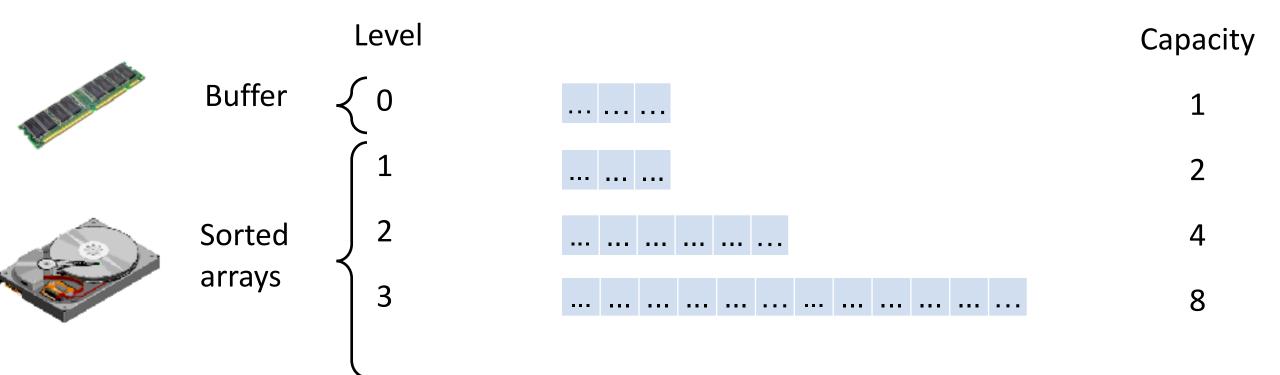
 $O(log_2 N/B)$ 

**Total search cost:** 

 $O(log_2(N/P) * log_2(N/B))$ 

If P≈B then:

 $O(\log_2(N/B)^2)$ 

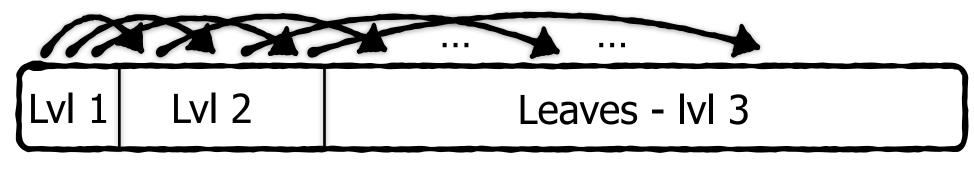


## Basic LSM-tree - Get Queries Cost

We can do slightly better by structuring each file (SST) as a static B-tree. How would this impact search cost?

Total search cost:

 $O(log_2(N/P) * log_2(N/B))$ 



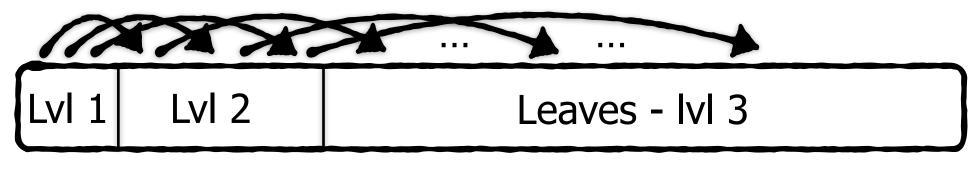
SST static B-tree structure

## Basic LSM-tree - Get Queries Cost

We can do slightly better by structuring each file (SST) as a static B-tree. How would this impact search cost?

Total search cost:

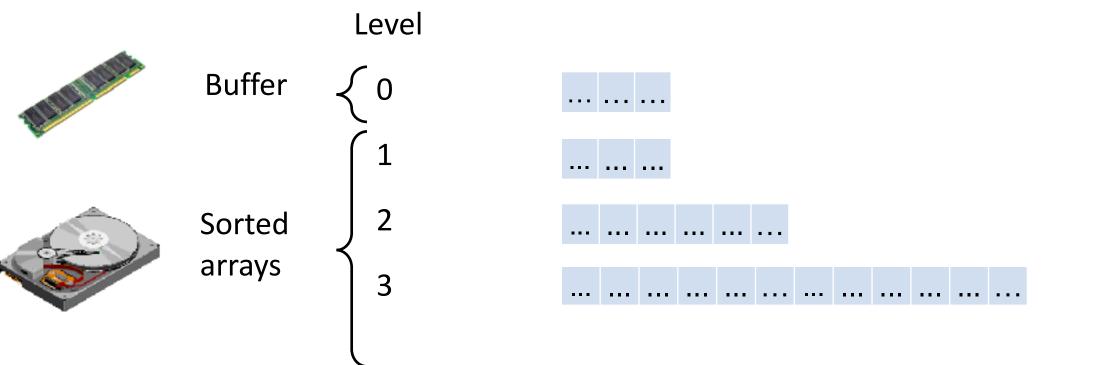
 $O(log_2(N/P) * log_B N)$ 



SST static B-tree structure

## Basic LSM-tree - Scan Queries

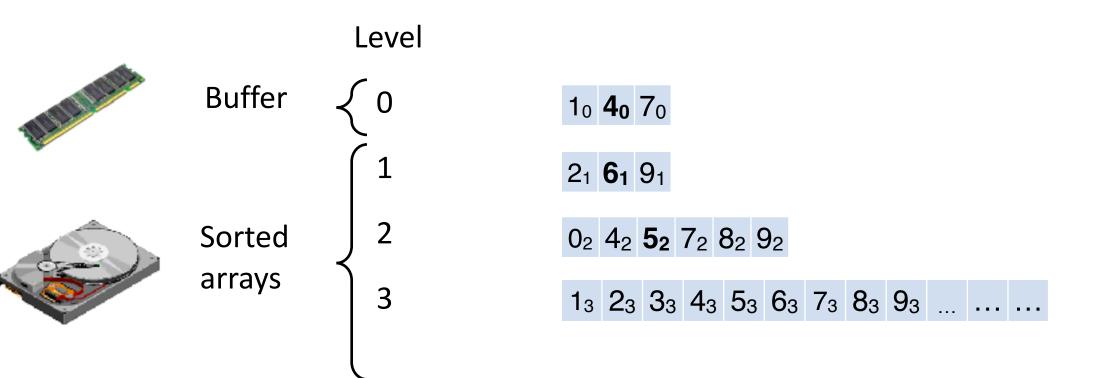
Return most recent version of each entry in the range across entire tree.



## Basic LSM-tree - Scan Queries

Return most recent version of each entry in the range across entire tree.

e.g., Scan(4, 6) - range inclusive Expected output?

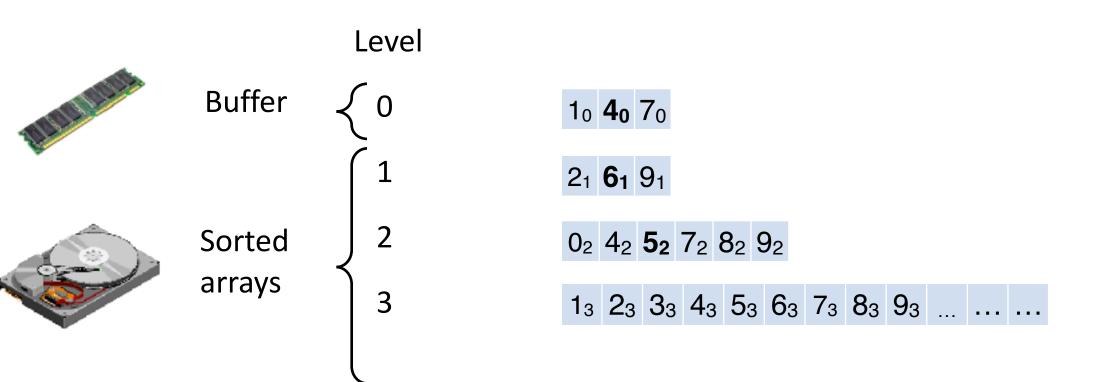


## Basic LSM-tree - Scan Queries

Return most recent version of each entry in the range across entire tree.

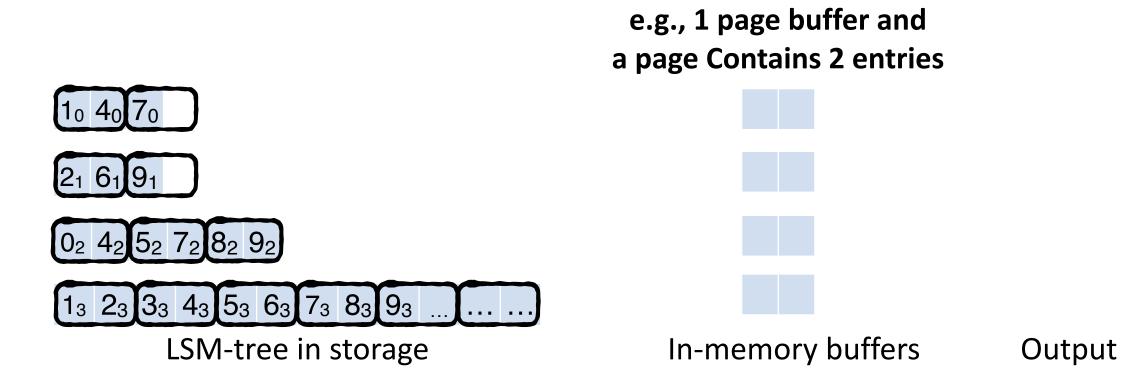
e.g., Scan(4, 6)

Expected output: 40 52 61



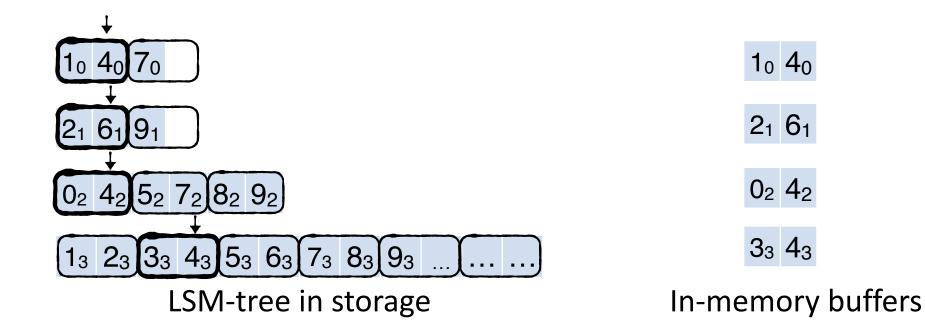
## Scan(4, 6)

1. Allocate an in-memory buffer (≥1 page) for each level



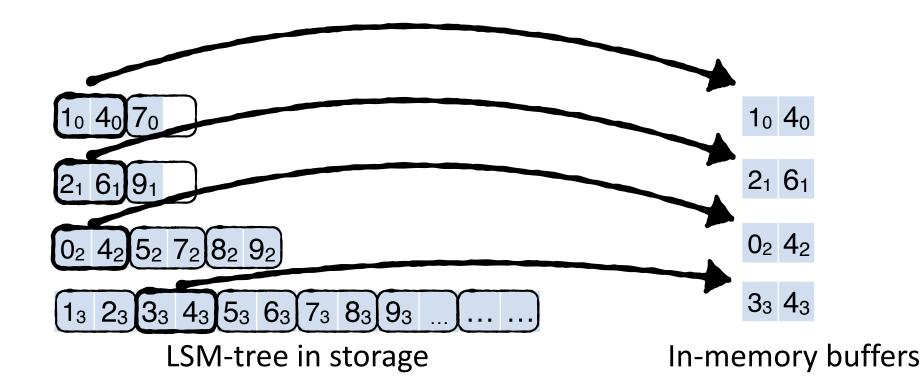
## Scan(4, 6)

- 1. Allocate an in-memory buffer (≥1 page) for each level
- 2. Search for start of key range at each level



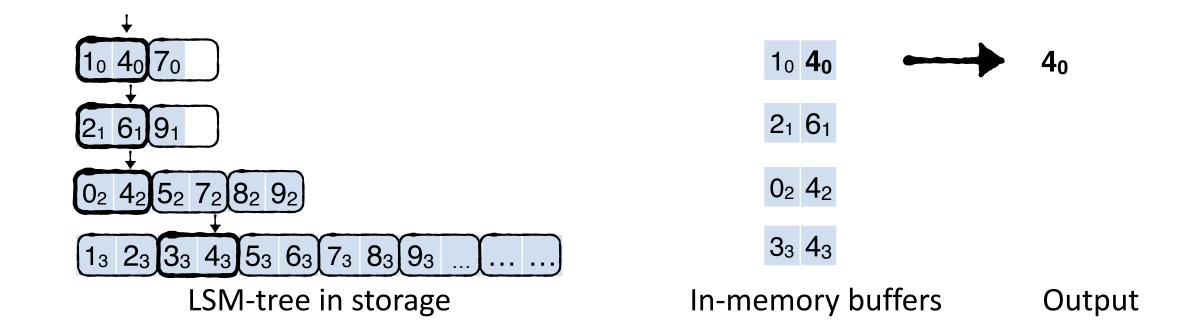
Output

- 1. Allocate an in-memory buffer (≥1 page) for each level
- 2. Search for start of key range at each level

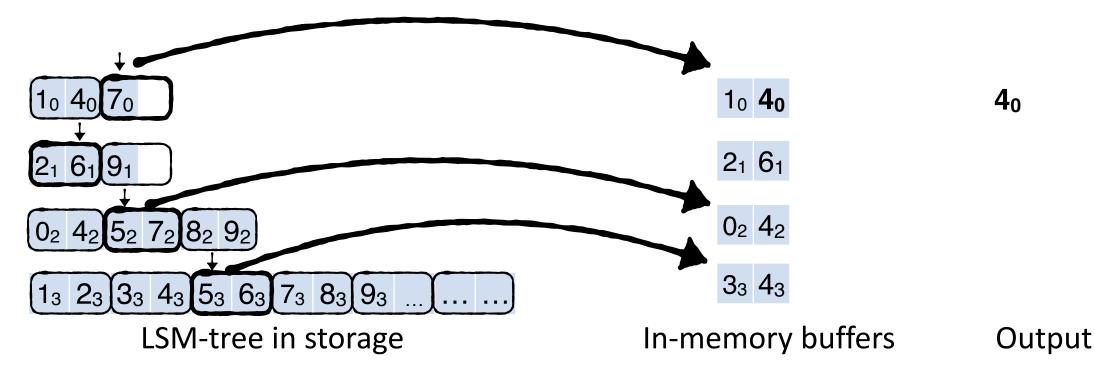


Output

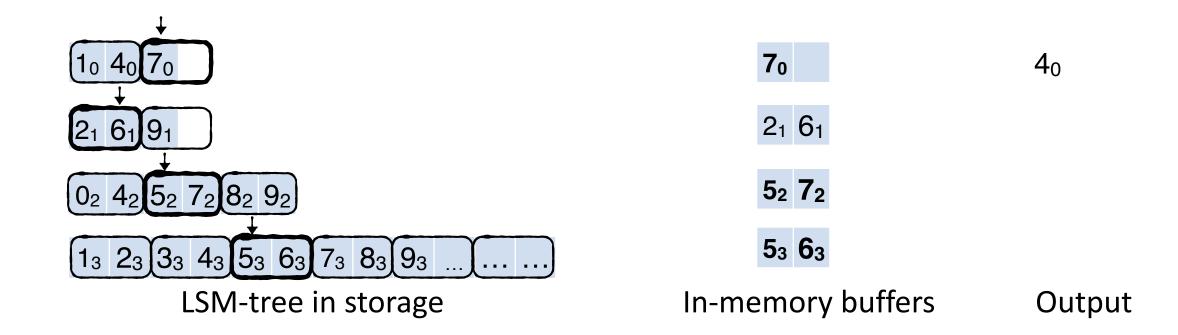
- 1. Allocate an in-memory buffer (≥1 page) for each level
- 2. Search for start of key range at each level Loop until reaching end of range
  - 3. Output youngest version of entry with next smallest key



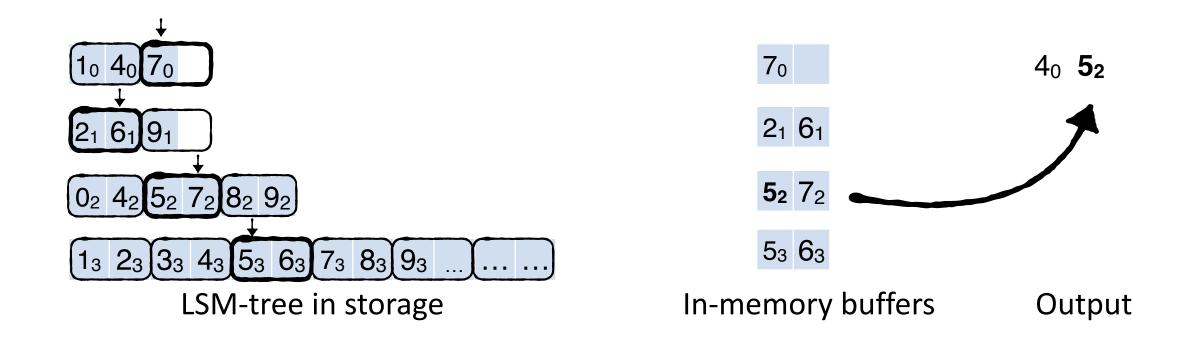
- 1. Allocate an in-memory buffer (≥1 page) for each level
- 2. Search for start of key range at each level Loop until reaching end of range
  - 3. Output youngest version of entry with next smallest key
  - 4. If we traverse last entry in a given buffer, read next page from run



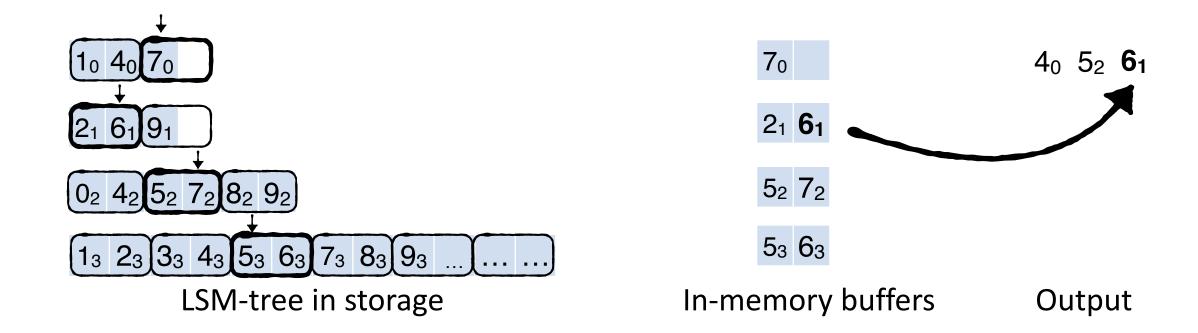
- 1. Allocate an in-memory buffer (≥1 page) for each level
- 2. Search for start of key range at each level Loop until reaching end of range
  - 3. Output youngest version of entry with next smallest key
  - 4. If we traverse last entry in a given buffer, read next page from run

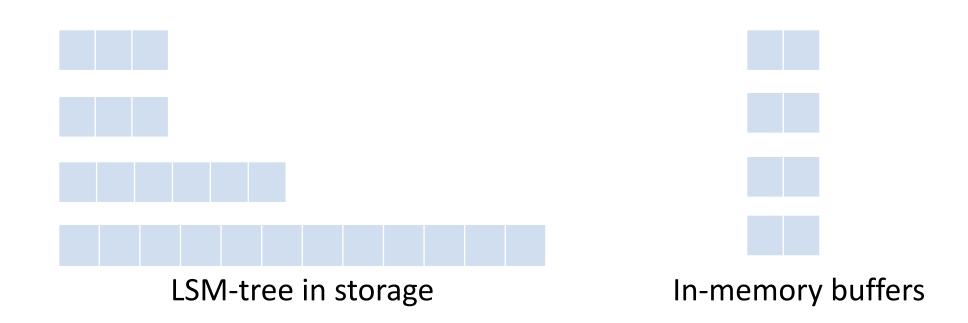


- 1. Allocate an in-memory buffer (≥1 page) for each level
- 2. Search for start of key range at each level Loop until reaching end of range
  - 3. Output youngest version of entry with next smallest key
  - 4. If we traverse last entry in a given buffer, read next page from run

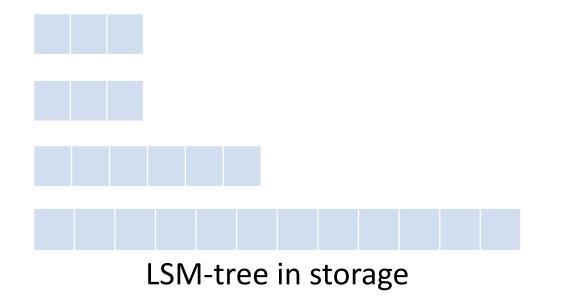


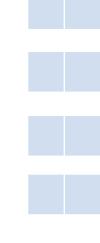
- 1. Allocate an in-memory buffer (≥1 page) for each level
- 2. Search for start of key range at each level Loop until reaching end of range
  - 3. Output youngest version of entry with next smallest key
  - 4. If we traverse last entry in a given buffer, read next page from run



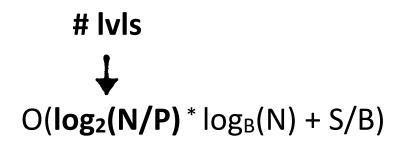


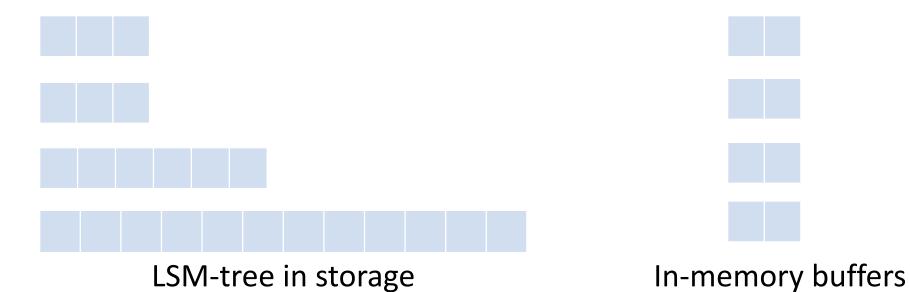
$$O(log_2(N/P) * log_B(N) + S/B)$$

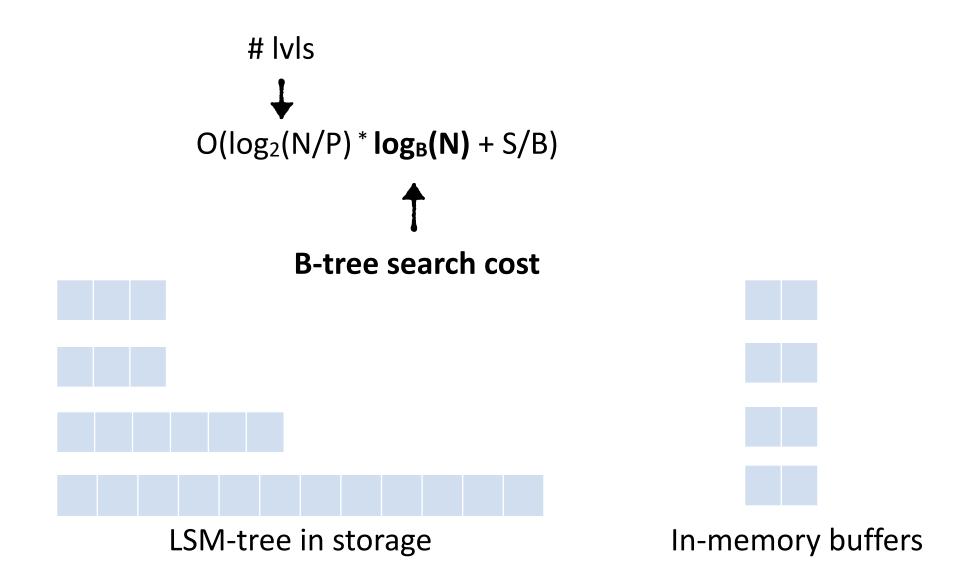


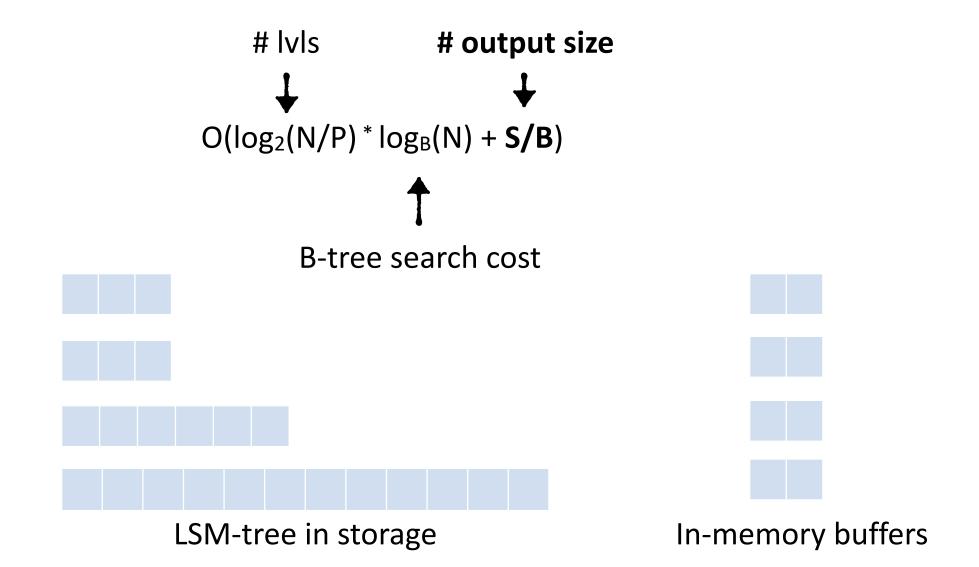


In-memory buffers

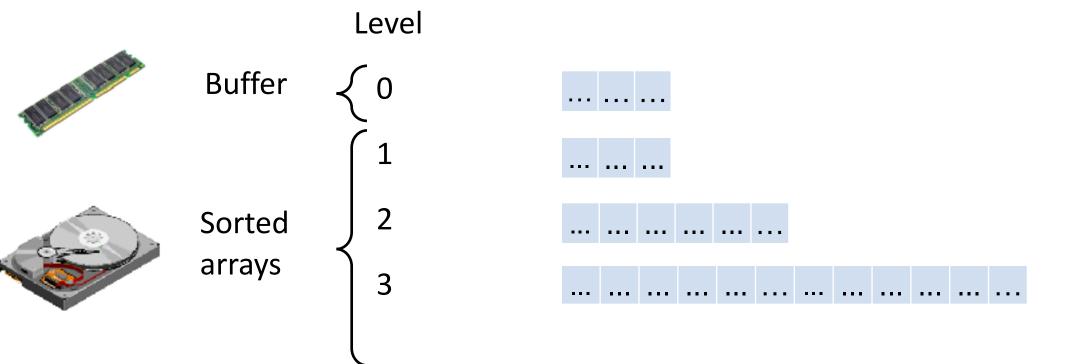




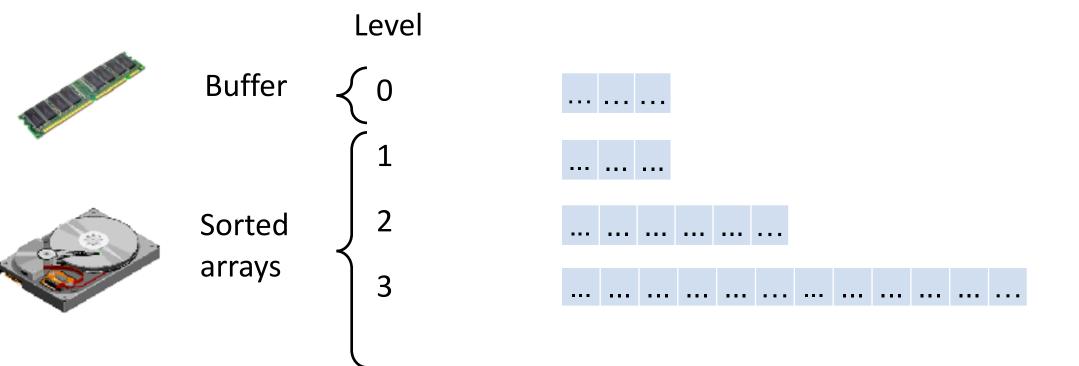




How many times is each entry copied?

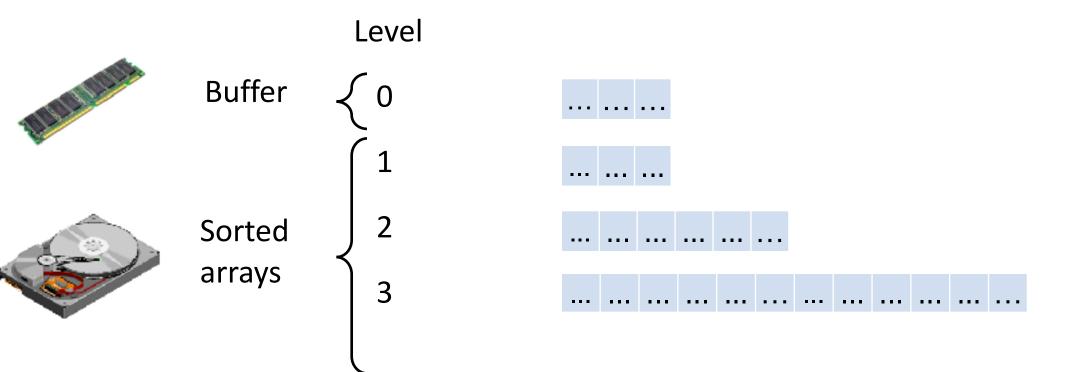


How many times is each entry copied?  $O(log_2 N/P)$ 



How many times is each entry copied?  $O(log_2 N/P)$ 

Price of each copy?

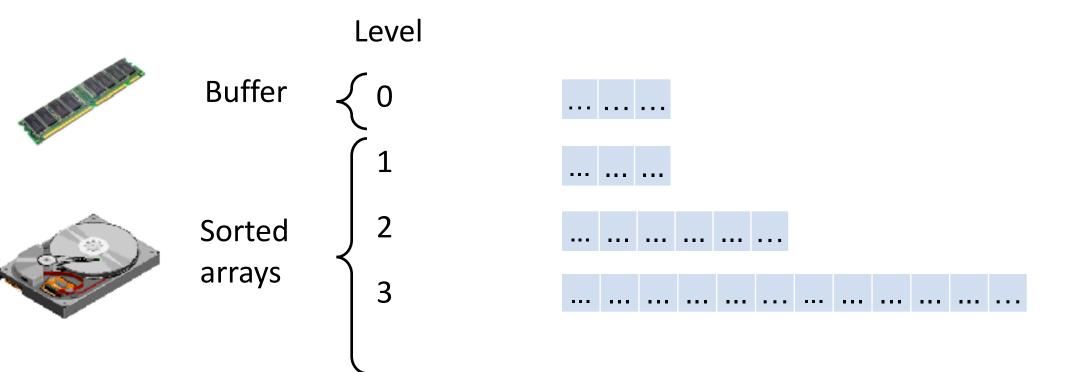


How many times is each entry copied?

 $O(log_2 N/P)$ 

Price of each copy?

O(1/B) reads & writes



How many times is each entry copied?

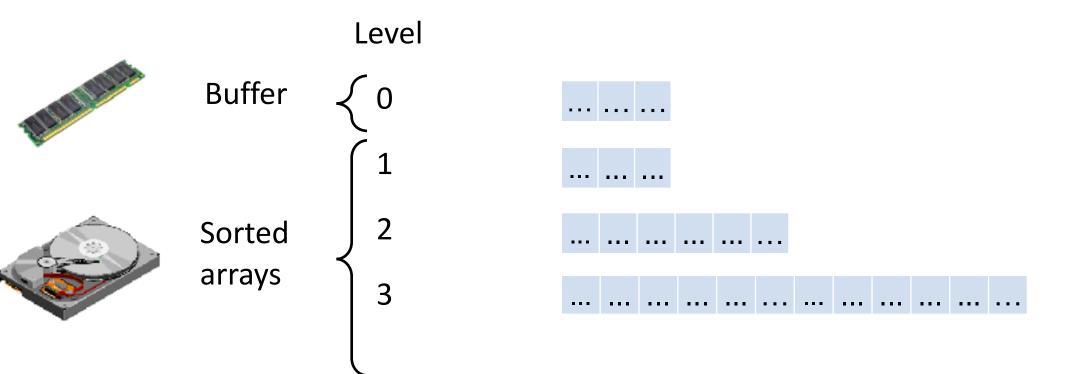
 $O(log_2 N/P)$ 

Price of each copy?

O(1/B) reads & writes

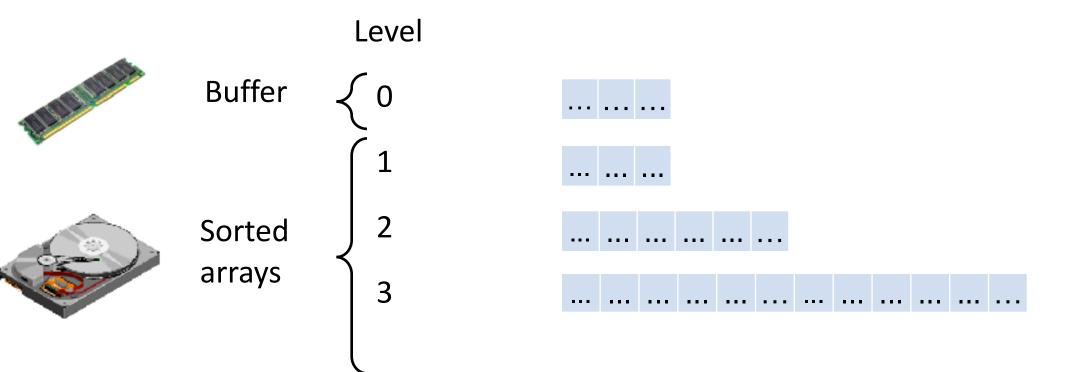
**Total cost:** 

O((log<sub>2</sub> N/P)/B) read & write I/Os



Total cost: O((log<sub>2</sub> N/P)/B) read & write I/Os

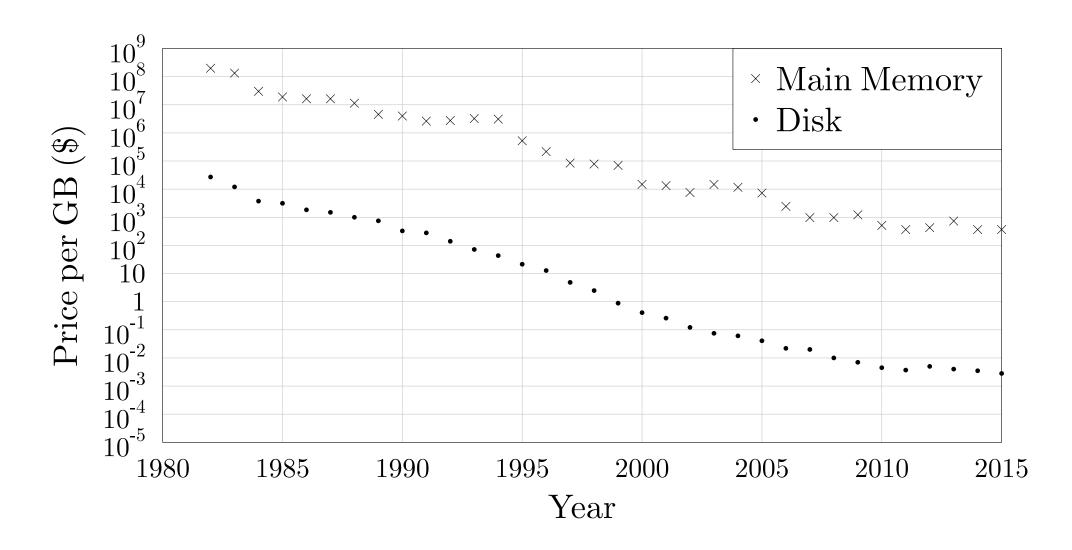
As all writes are large & sequential (rather than random), there is less SSD garbage-collection



Operation	I/O	append-only table	Basic LSM-tree table	B-Tree
Query	Reads	O(N/B)	O(log <sub>2</sub> (N/P) * log <sub>B</sub> N)	O(log <sub>B</sub> N)
Insert	Reads	0	O((log <sub>2</sub> N/P)/B)	O(log <sub>B</sub> N)
	Writes	O(1/B)	O((log <sub>2</sub> N/P)/B)	O(1) & GC

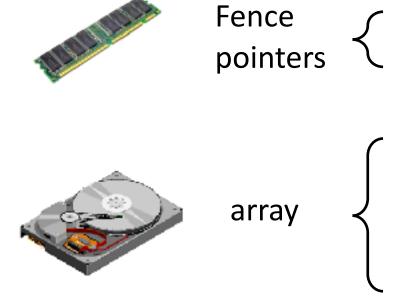
Break:)

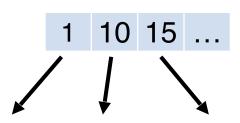
# Declining Main Memory Cost



# **Declining Main Memory Cost**

It's viable to pin the internal nodes of B-trees in memory





Block 1	Block 2	Block 3	
1	10	15	
3	11	16	•••
6	13	18	•••

Cost Metric	Unit	append-only table	Basic LSM-tree table	B-Tree
Query	Reads IOs	O(N/B)	O( log <sub>2</sub> (N/P) ************************************	O(logist)
Insert	Reads IOs	0	$O((log_2 N/P)/B)$	
	Write IOs	O(1/B)	$O((log_2 N/P)/B)$	O(1) & GC

Cost Metric	Unit	append-only table	Basic LSM-tree table	B-Tree
Query	Reads IOs	O(N/B)	O( log <sub>2</sub> (N/P) )	O(1)
Insert	Reads IOs	0	O((log <sub>2</sub> N/P)/B)	O(1)
	Write IOs	O(1/B)	O((log <sub>2</sub> N/P)/B)	O(1) & GC
Memory	#entries	O(B)	O(N/B)	O(N/B)

Memory cost is measured here as # of entries stored in memory

Basic LSM-tree

Tiered LSM-tree

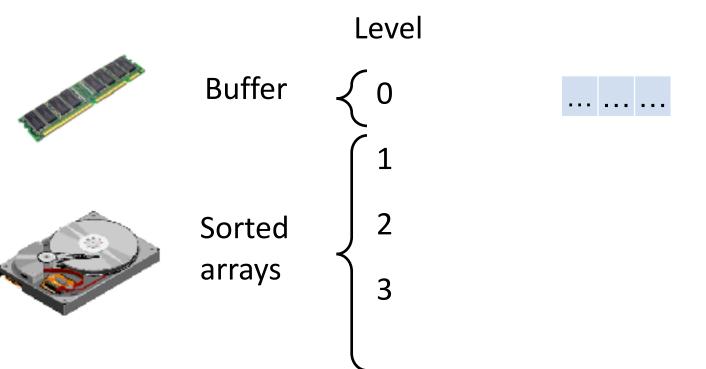




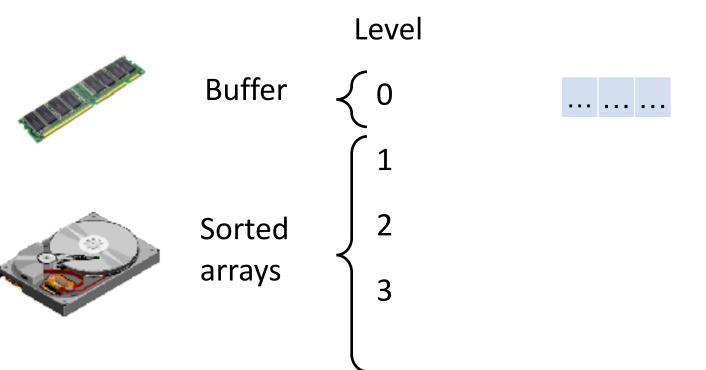




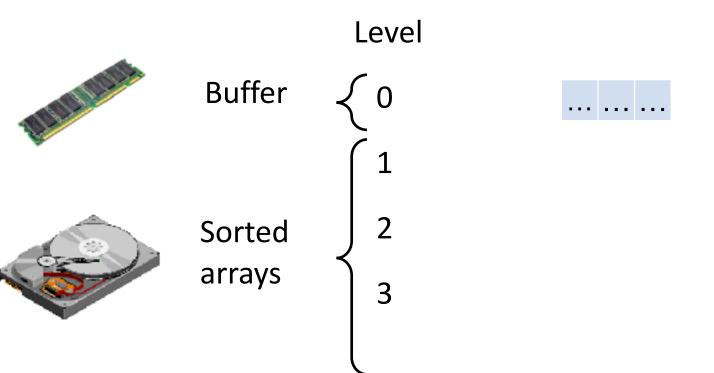
Lookup cost depends on number of levels



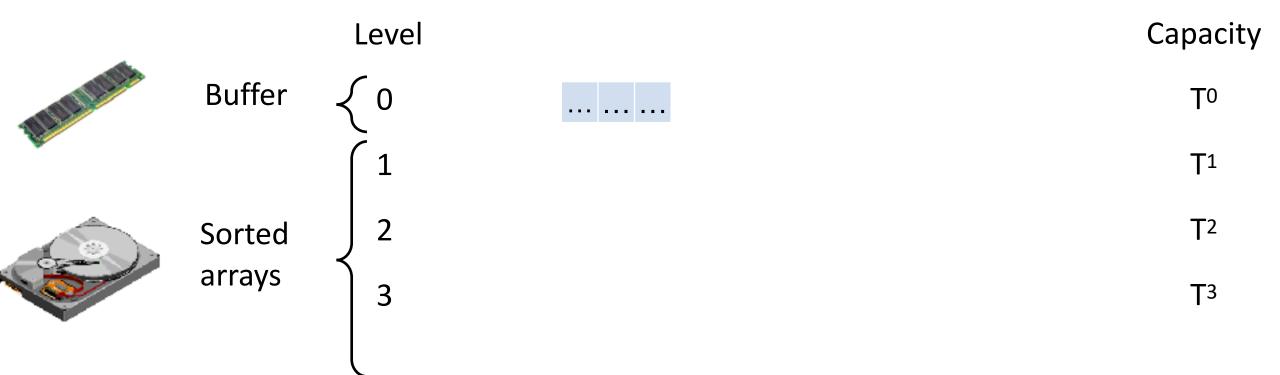
Lookup cost depends on number of levels How to reduce it?



Lookup cost depends on number of levels How to reduce it?



Lookup cost depends on number of levels How to reduce it?



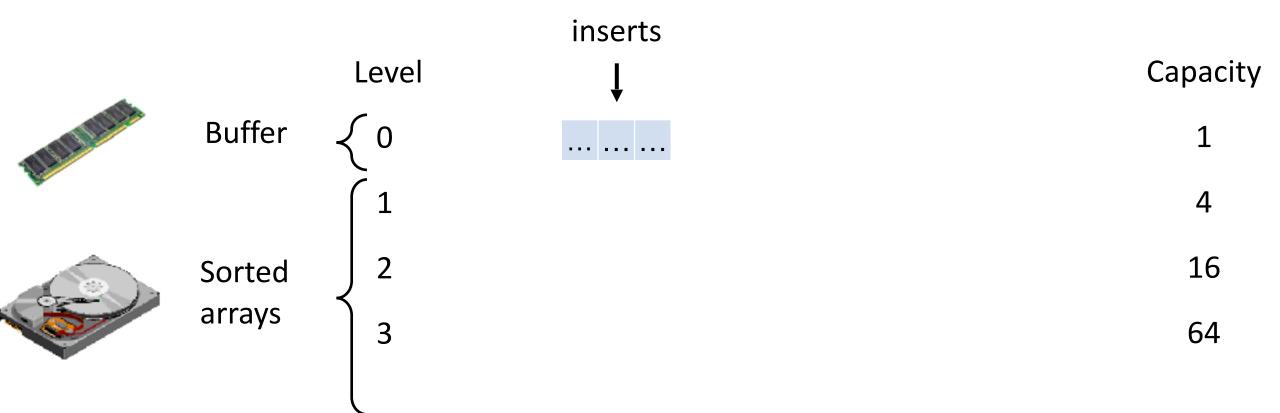
Lookup cost depends on number of levels How to reduce it?

E.g. size ratio of 4



Lookup cost depends on number of levels How to reduce it?

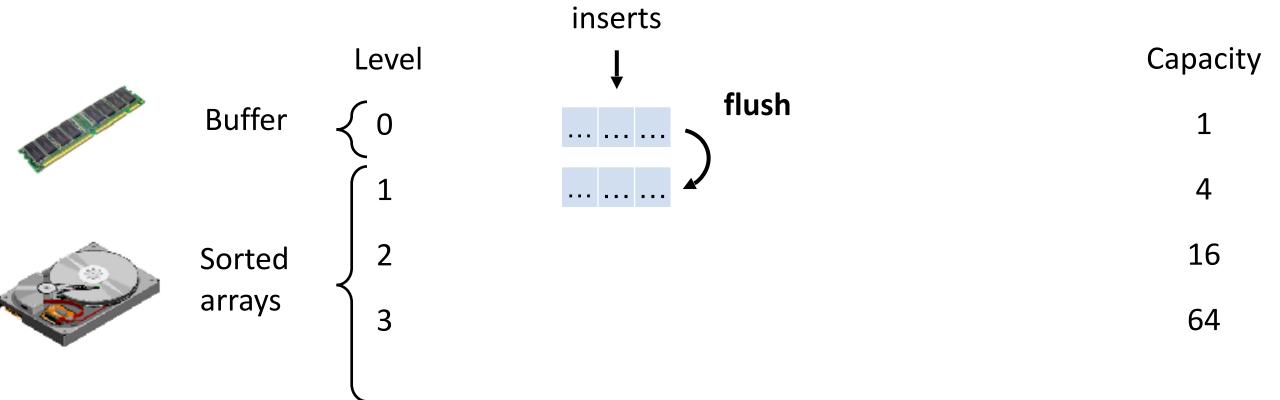
E.g. size ratio of 4



Lookup cost depends on number of levels

How to reduce it?

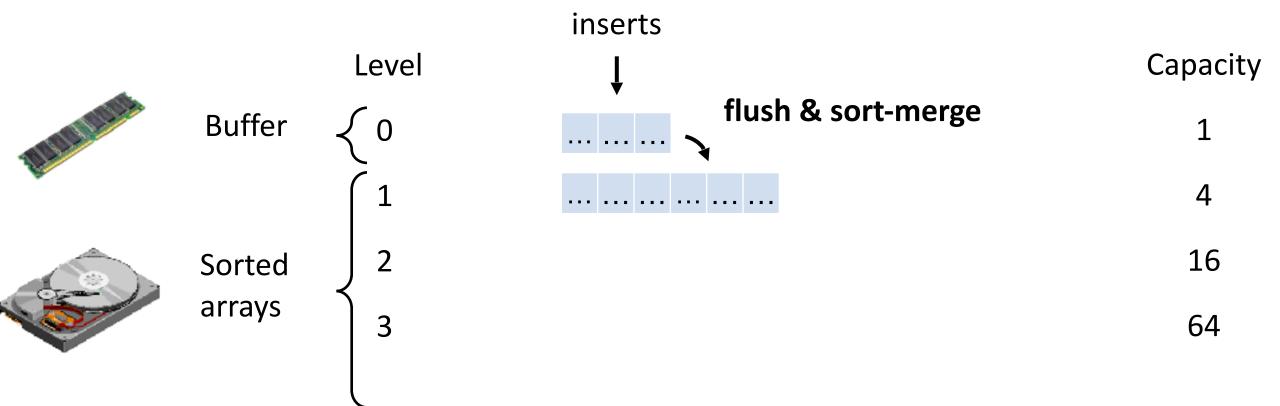
E.g. size ratio of 4



Lookup cost depends on number of levels

How to reduce it?

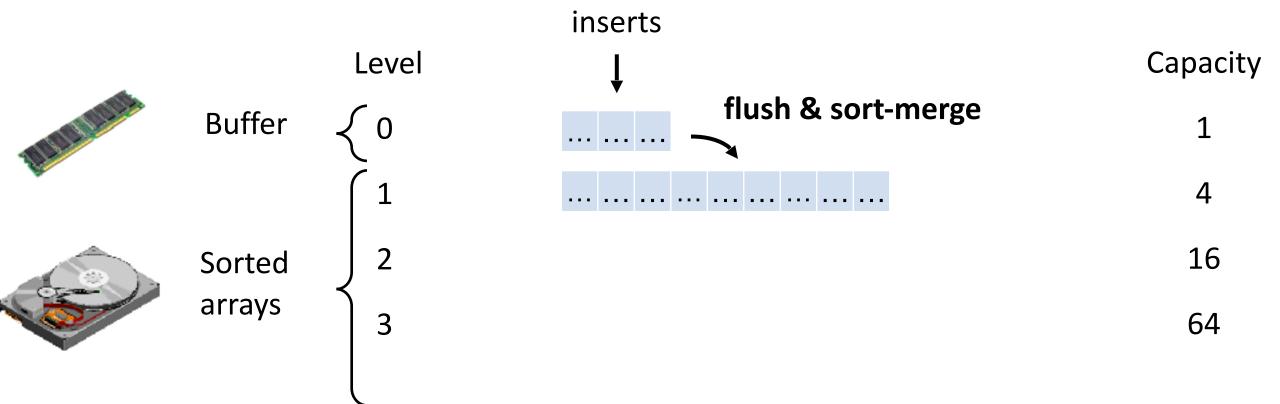
E.g. size ratio of 4



Lookup cost depends on number of levels

How to reduce it?

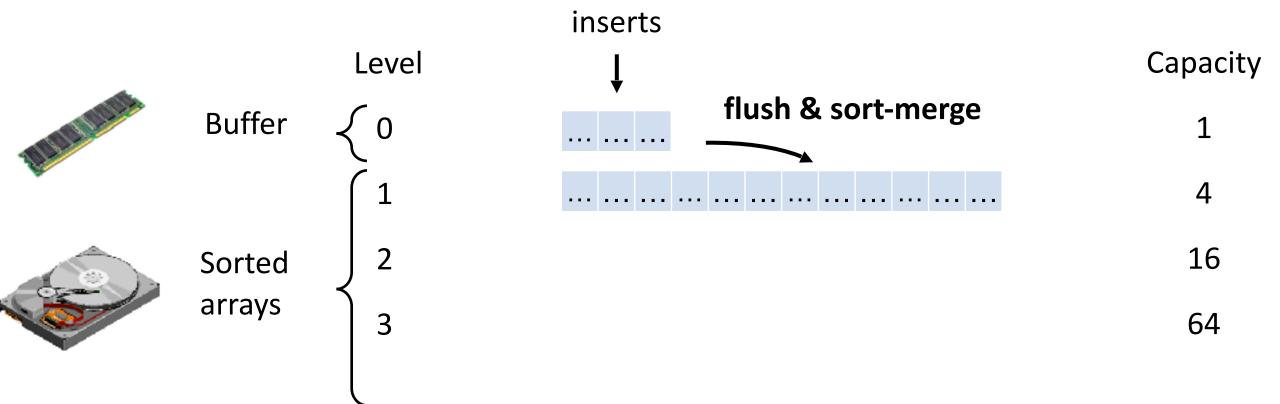
E.g. size ratio of 4



Lookup cost depends on number of levels

How to reduce it?

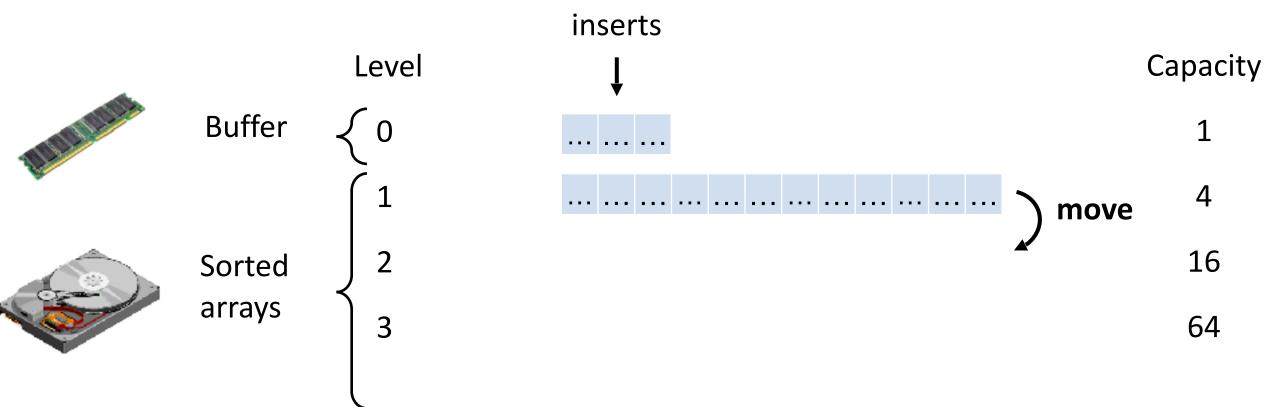
E.g. size ratio of 4



Lookup cost depends on number of levels

How to reduce it?

E.g. size ratio of 4

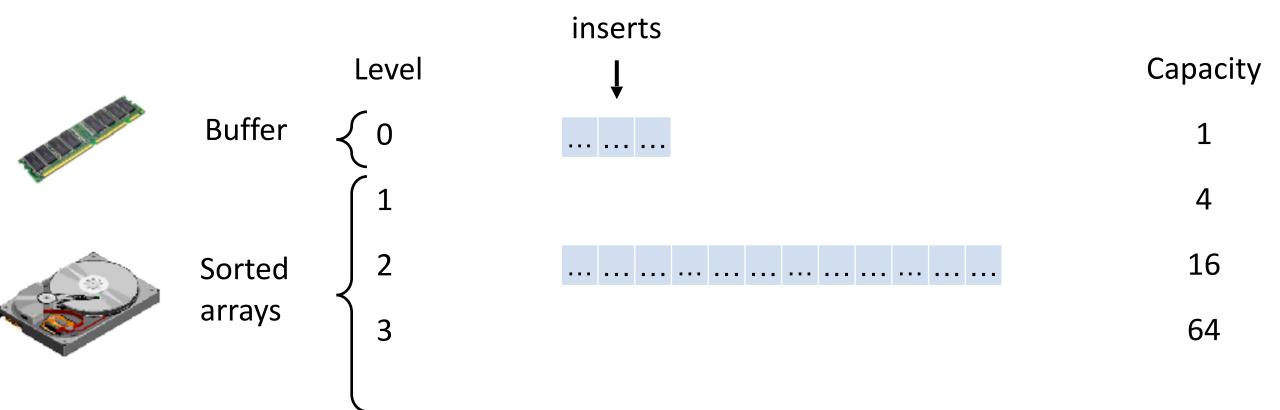


Lookup cost depends on number of levels

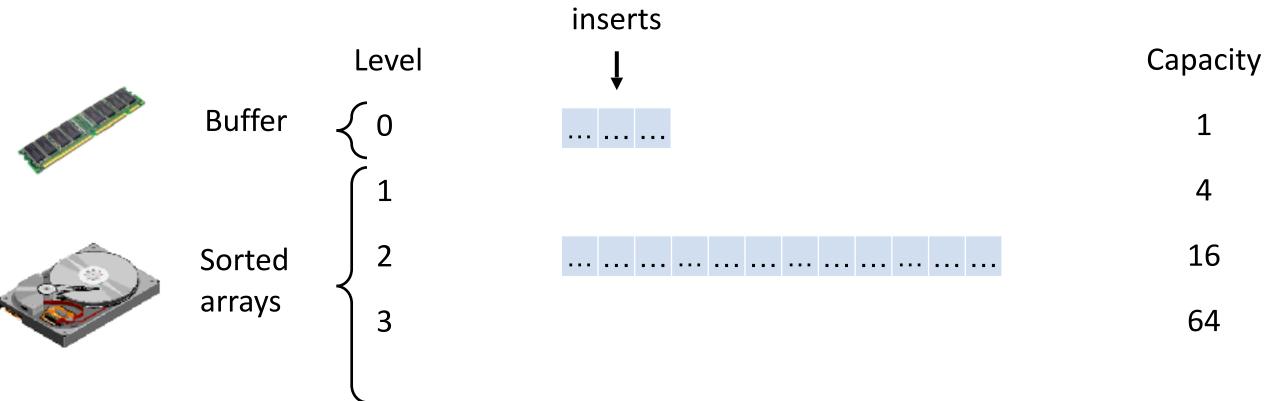
How to reduce it?

E.g. size ratio of 4

Increase size ratio T

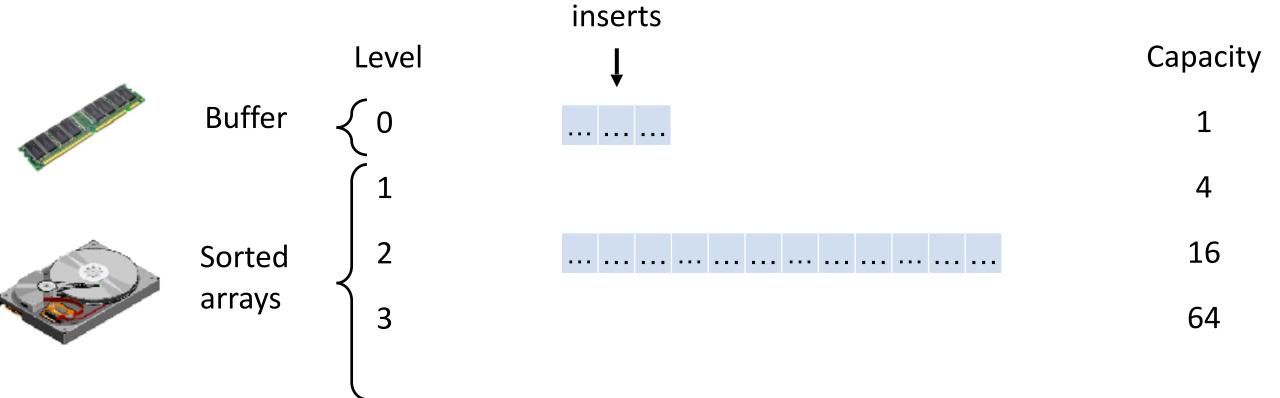


Lookup cost?



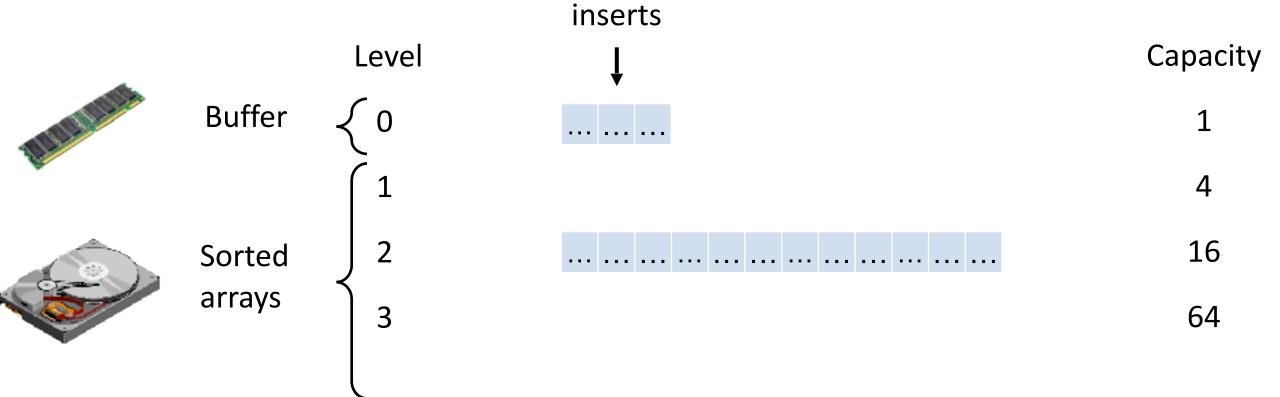
Lookup cost?

 $O(log_T(N/P))$ 



Lookup cost?
O(log<sub>T</sub>(N/P))

Insertion cost?

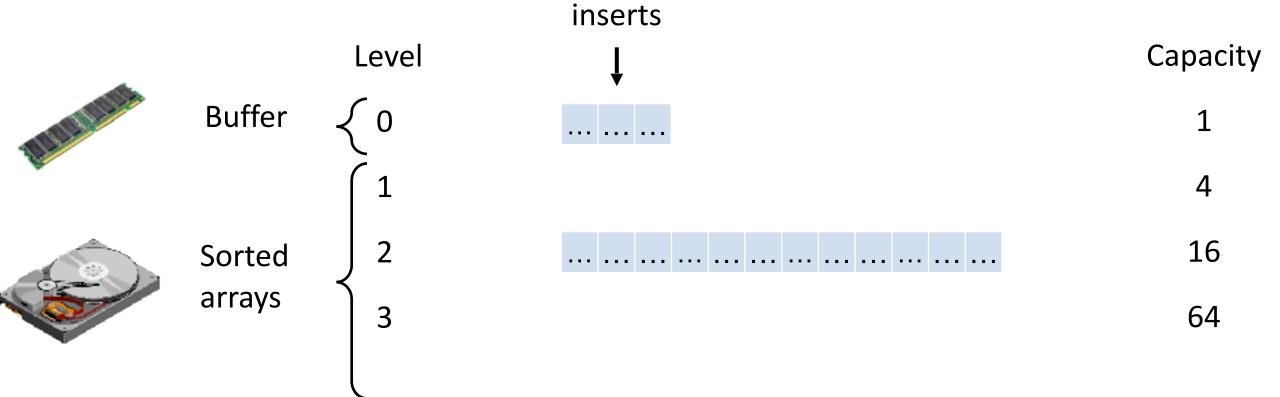


Lookup cost?

 $O(\log_T(N/P))$ 

Insertion cost?

 $O(T/B \bullet log_T(N/P))$ 



Lookup cost?

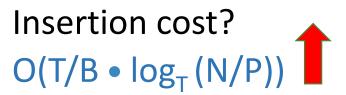
 $O(\log_T(N/P))$ 

Insertion cost?

 $O(T/B \bullet log_T(N/P))$ 

What happens as we increase the size ratio T?

Lookup cost?
O(log<sub>T</sub>(N/P))



What happens as we increase the size ratio T?

Lookup cost?
O(log<sub>T</sub>(N/P))

Insertion cost?
O(T/B • log<sub>T</sub> (N/P))

What happens as we increase the size ratio T?

What happens when size ratio T is set to be N/P?



Lookup cost?  $O(log_T(N/P))$ 

Insertion cost?
O(T/B • log<sub>T</sub> (N/P))

What happens as we increase the size ratio T?

What happens when size ratio T is set to be N/P?

Lookup cost becomes:

O(1)

Insert cost becomes:

 $O(N/(B \cdot P))$ 



Lookup cost?
O(log<sub>T</sub>(N/P))

Insertion cost?
O(T/B • log<sub>T</sub> (N/P))

What happens as we increase the size ratio T?

What happens when size ratio T is set to be N/P?

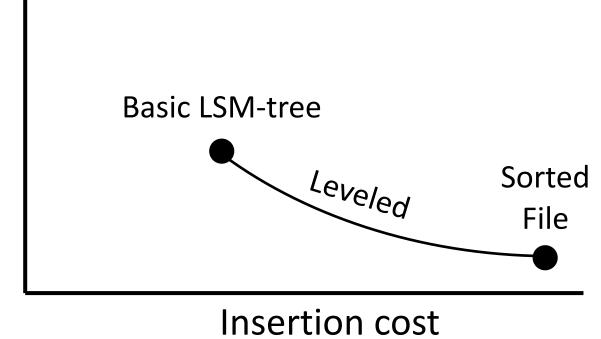
Lookup cost becomes:

O(1)

Insert cost becomes:

 $O(N/(B \cdot P))$ 

The LSM-tree becomes a sorted file!



Basic LSM-tree

Tiered LSM-tree









Reduce the number of levels by increasing the size ratio.

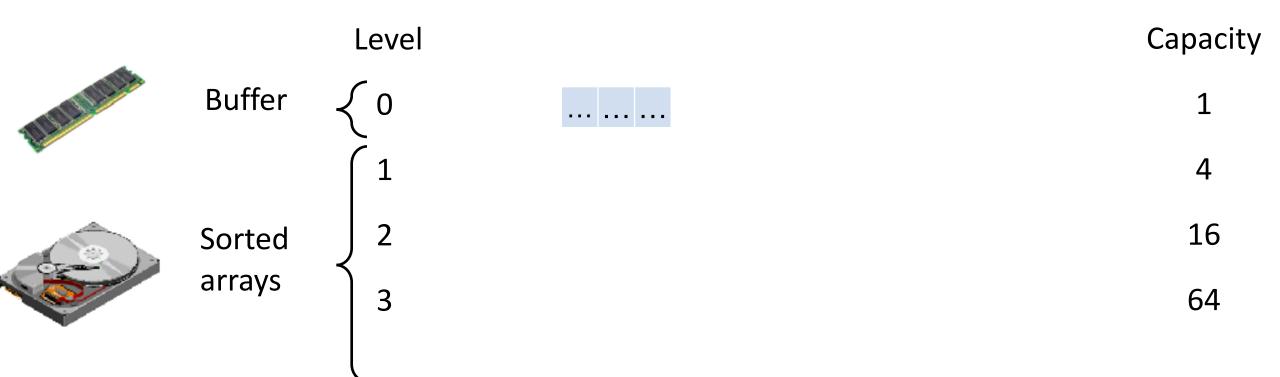


Reduce the number of levels by increasing the size ratio. Do not merge within a level.



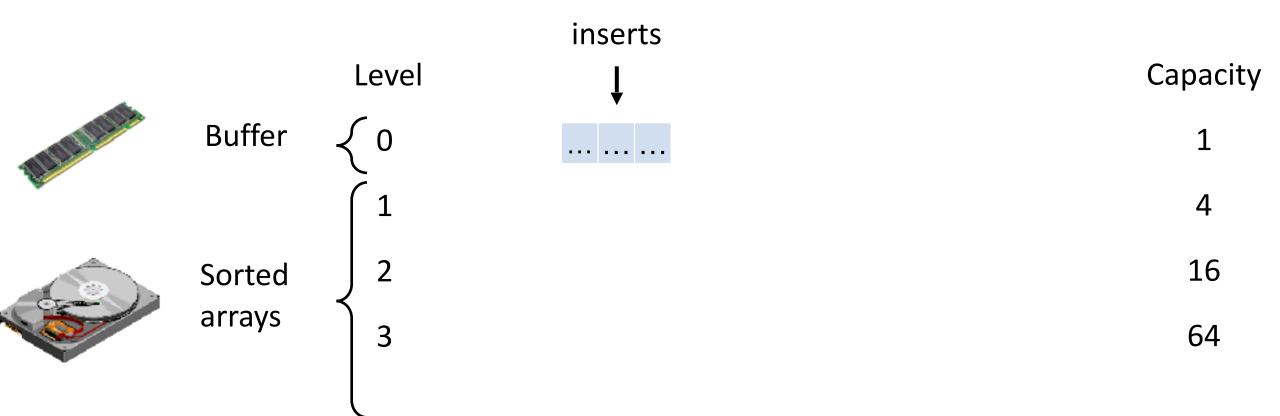
Reduce the number of levels by increasing the size ratio.

Do not merge within a level.



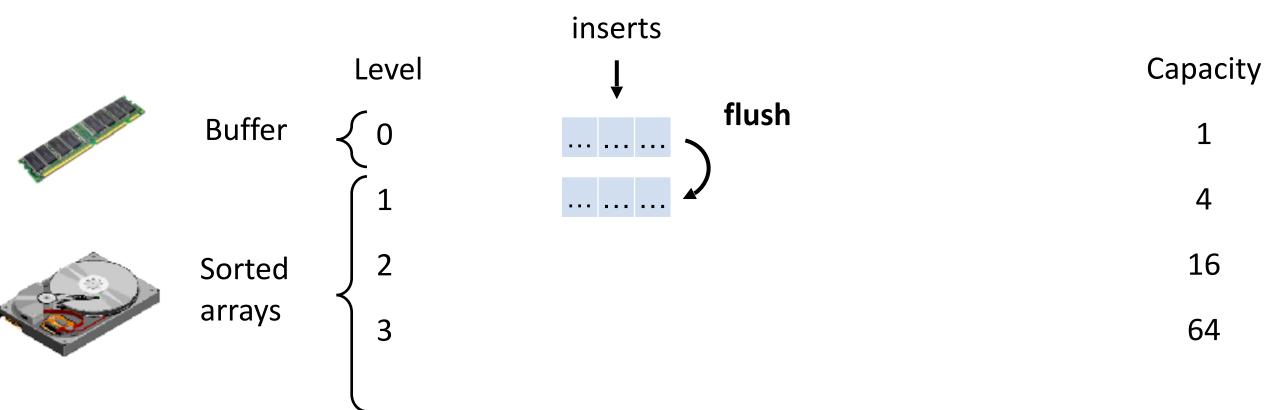
Reduce the number of levels by increasing the size ratio.

Do not merge within a level.



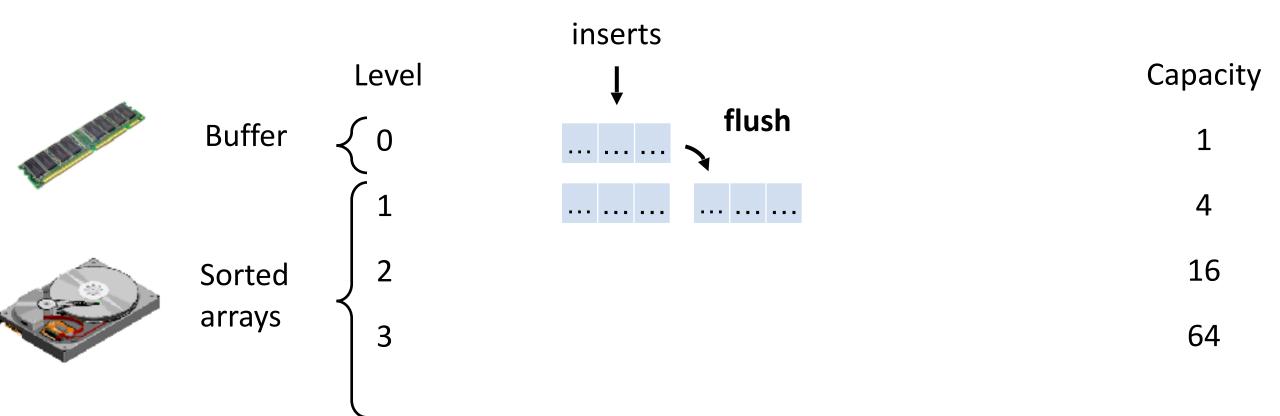
Reduce the number of levels by increasing the size ratio.

Do not merge within a level.



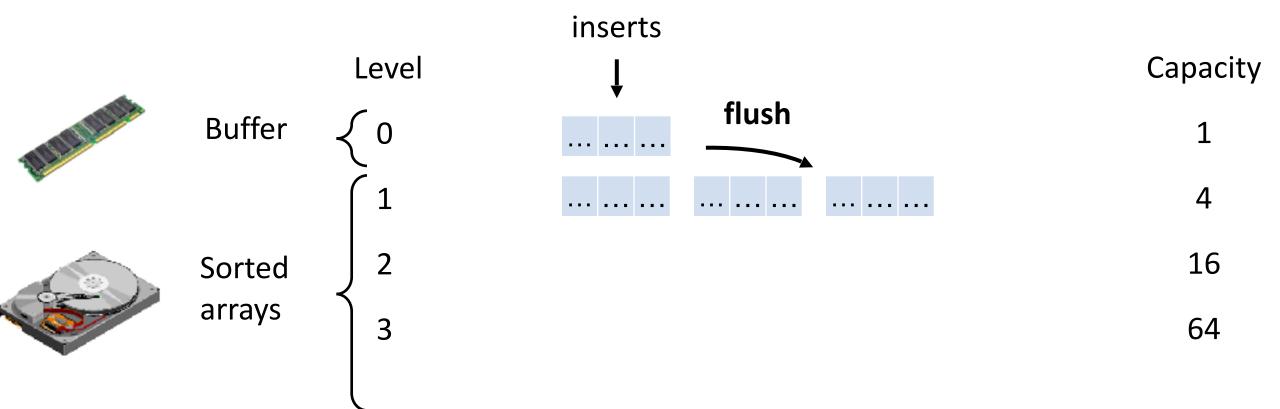
Reduce the number of levels by increasing the size ratio.

Do not merge within a level.



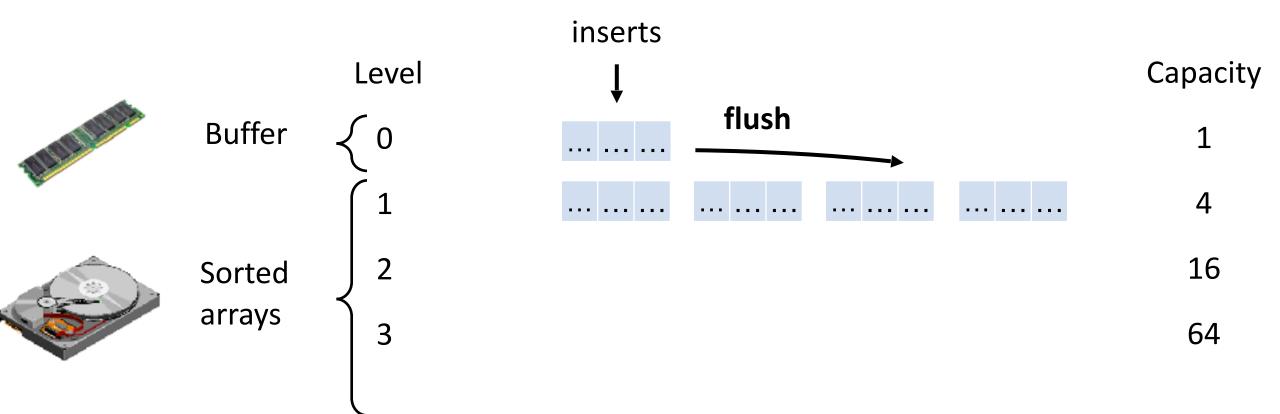
Reduce the number of levels by increasing the size ratio.

Do not merge within a level.



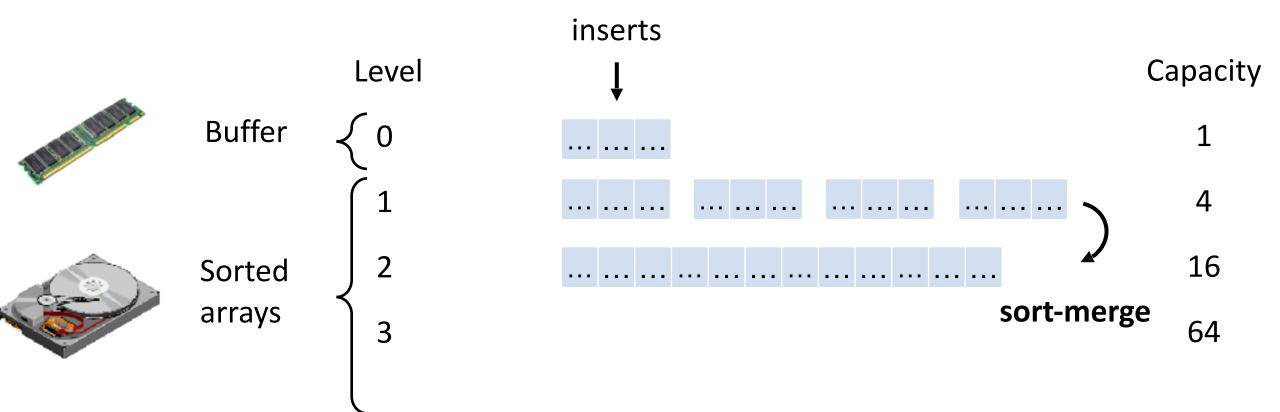
Reduce the number of levels by increasing the size ratio.

Do not merge within a level.



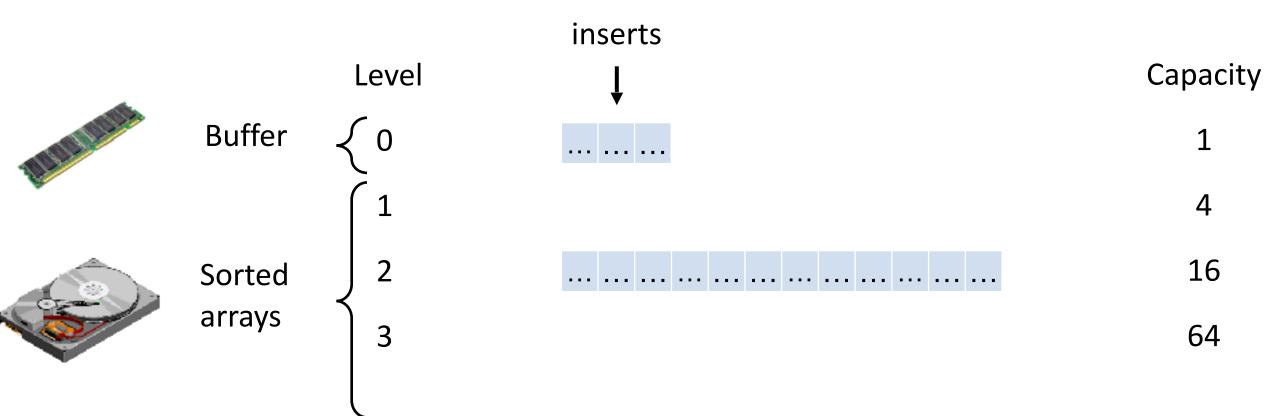
Reduce the number of levels by increasing the size ratio.

Do not merge within a level.

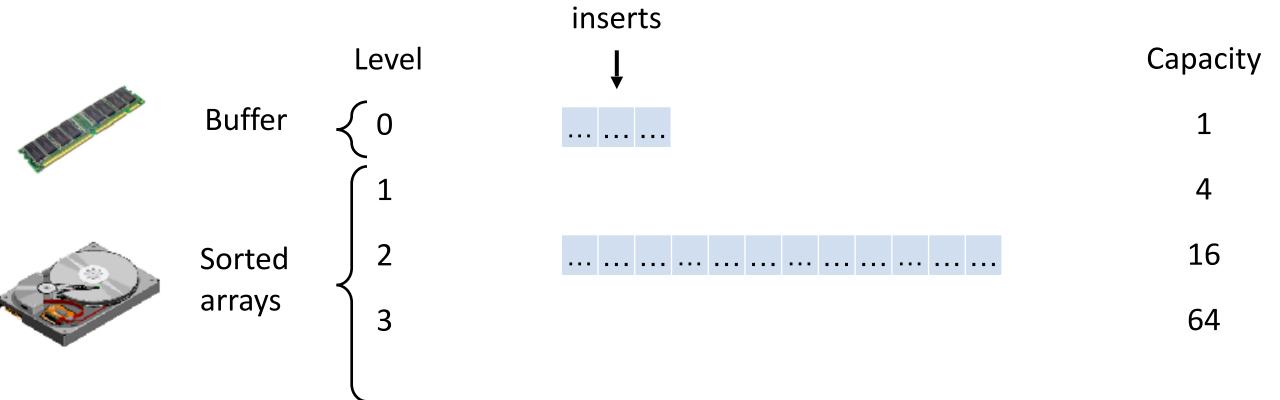


Reduce the number of levels by increasing the size ratio.

Do not merge within a level.

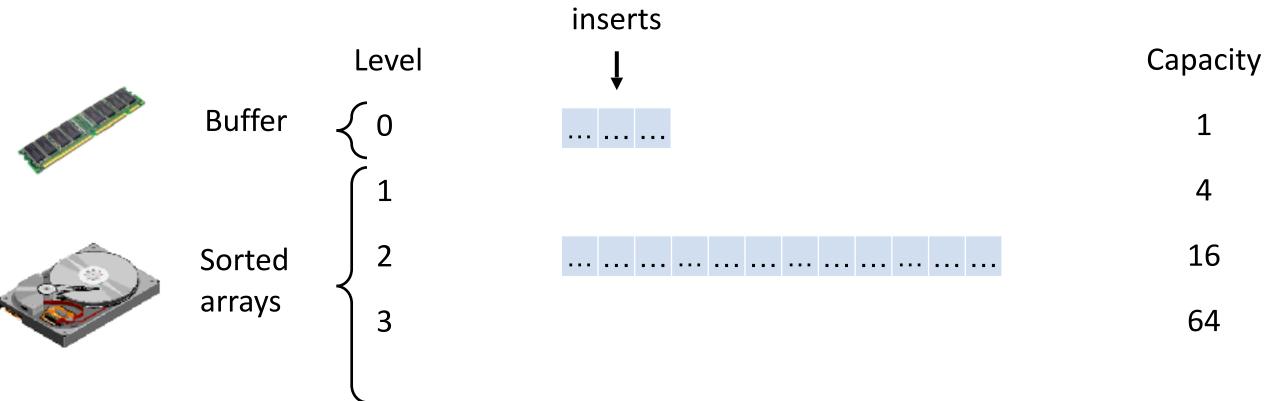


Lookup cost?



Lookup cost?

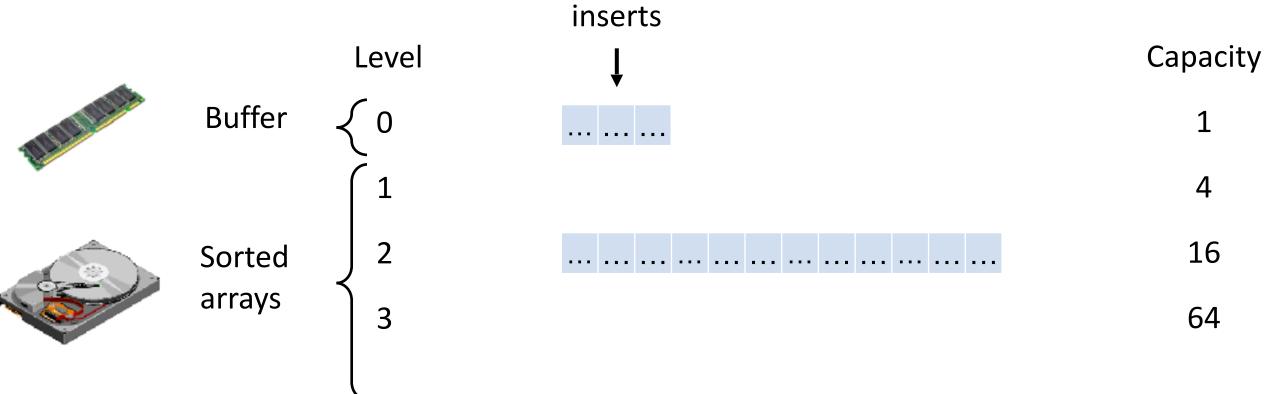
 $O(T \bullet log_T(N/P))$ 



Lookup cost?

 $O(T \bullet log_T(N/P))$ 

Insertion cost?

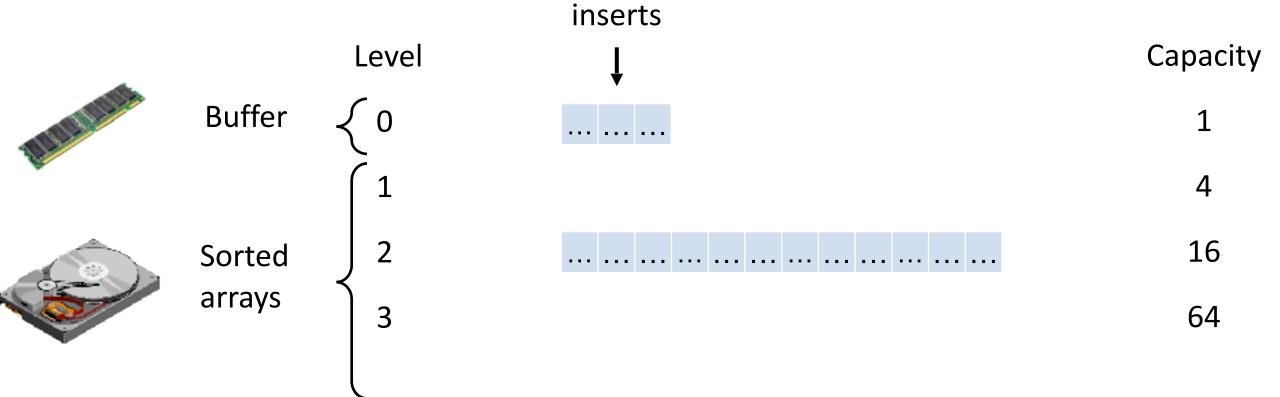


Lookup cost?

 $O(T \bullet log_T(N/P))$ 

Insertion cost?

 $O(1/B \bullet log_T(N/P))$ 



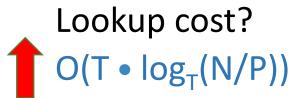
Lookup cost?

 $O(T \bullet log_T(N/P))$ 

Insertion cost?

 $O(1/B \bullet log_T(N/P))$ 

What happens as we increase the size ratio T?



Insertion cost?

 $O(1/B \bullet log_T(N/P))$ 



What happens as we increase the size ratio T?



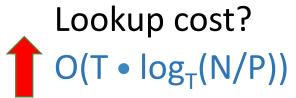
Insertion cost?

 $O(1/B \bullet log_T(N/P))$ 



What happens as we increase the size ratio T?

What happens when size ratio T is set to be N/P?



Insertion cost?

 $O(1/B \bullet log_T(N/P))$ 



What happens as we increase the size ratio T?

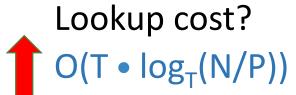
What happens when size ratio T is set to be N/P?

Lookup cost becomes:

O(N/P)

Insert cost becomes:

O(1/B)



Insertion cost?

 $O(1/B \bullet log_T(N/P))$ 



What happens as we increase the size ratio T?

What happens when size ratio T is set to be N/P?

Lookup cost becomes:

O(N/P)

Insert cost becomes:

O(1/B)

The tiered LSM-tree becomes an append-only file!

Insertion cost

Cost Metric	Unit	Unordered File	Tiered LSM-tree	Leveled LSM-tree	B-Tree
Query	Reads IOs	O(N/B)	O(L * T)	O(L)	O(1)
Insert	Reads IOs	0	O(L/B)	O((L * T)/B)	O(1)
	Write IOs	O(1/B)	O(L/B)	O((L * T)/B)	O(1) & GC
Memory	#entries	O(B)	O(N/B)	O(N/B)	O(N/B)

Let  $L = log_T(N/P)$ 

This table assumes internal nodes fit in memory.

Write-optimized

Write-optimized

Highly tunable

Write-optimized

Highly tunable

Backbone of many modern systems

Write-optimized

Highly tunable

Backbone of many modern systems

Trade-off between lookup and insert cost (tiering/leveling, size ratio)

Write-optimized

Highly tunable

Backbone of many modern systems

Trade-off between lookup and insert cost (tiering/leveling, size ratio)

Trade main memory for lookup cost (fence pointers, Bloom filters)