

# HashMap Internal Structure



How does a HashMap work internally?  
Discuss the underlying data structures and the concept of hashing.



A `HashMap` in Java (or similar hash-based data structures in other programming languages) is a key-value pair-based collection that allows fast retrieval, insertion, and deletion of elements. Internally, it works by leveraging the concept of **hashing**, which maps keys to specific indices in an array. Here's a breakdown of how it works internally:

## 1. Underlying Data Structure:

- **Array of Buckets:** Internally, a `HashMap` uses an array where each element is referred to as a **bucket**. This array is where the key-value pairs (entries) are stored.
- Each bucket holds either a `null` (if no element has been placed in that bucket) or a reference to the entries stored in that bucket. Initially, the array size is often set to a default capacity, typically 16 in Java.
- **Entry Object:** Each entry in the `HashMap` is usually a node that stores the following:
  - Key
  - Value
  - A reference to the next node in case of hash collisions (i.e., when multiple keys hash to the same bucket).

## 2. Hashing:

- **Hash Function:** When a key-value pair is inserted into the `HashMap`, the key is passed through a hash function. This function computes an integer hash code that represents the key.
- The hash code is then **modulo-ed** with the size of the array to determine the index (bucket) where the key-value pair should be stored.
  - $\text{Index} = \text{hash}(\text{key}) \% \text{array.length}$
- In Java, the hash function used is derived from the `hashCode()` method of the key's class. However, the raw hash value is typically processed (e.g., via bitwise operations) to minimize collisions and ensure better distribution.

### 3. Collision Handling:

- When two different keys produce the same hash code (or when `hash(key) % array.length` returns the same index), a **collision** occurs. `HashMap` handles collisions using one of the following techniques:
  - **Separate Chaining (Linked List)**: Each bucket contains a linked list of entries. If multiple keys hash to the same bucket, they are linked together as nodes in this list. When retrieving a value, the list is traversed, and each entry's key is compared using the `equals()` method.
  - **Tree-Based Structure (Since Java 8)**: If the number of elements in a bucket's linked list exceeds a threshold (usually 8), the list is converted into a **balanced binary search tree** (red-black tree). This improves performance from  $O(n)$  to  $O(\log n)$  for lookups in the case of large numbers of collisions.

### 4. Operations in a HashMap:

- **Put Operation (Insertion)**:
  1. The key is hashed, and the index of the array (bucket) is determined.
  2. If the bucket is empty, the entry is placed directly there.
  3. If a collision occurs, the new entry is appended to the linked list or tree in that bucket.
  4. If the key already exists, its value is updated.
- **Get Operation (Retrieval)**:
  1. The key is hashed to find the corresponding bucket.
  2. The bucket is checked. If it's a linked list or a tree, it is traversed until the entry with the matching key is found.
  3. If found, the value is returned; otherwise, `null` is returned.
- **Remove Operation (Deletion)**:
  1. The key is hashed, and the index is determined.
  2. The corresponding bucket is searched for the entry with the given key.
  3. The entry is removed from the bucket (whether from a list or tree).

### 5. Resizing the HashMap:

- As more entries are added, the `HashMap` may need to grow to maintain efficient performance. This is done via **resizing** when the load factor (number of elements/size of the array) exceeds a predefined threshold (typically 0.75).
- During resizing, the `HashMap` creates a new array with double the previous capacity, and all existing entries are **rehashed** and redistributed into the new array.

### 6. Time Complexity of Operations:

- **Best Case:**  $O(1)$  for insertion, deletion, and retrieval if the hash function distributes the keys uniformly and there are few collisions.
- **Worst Case:**  $O(n)$  for insertion, deletion, and retrieval if all keys hash to the same bucket, creating a long linked list (or a tree in modern implementations, making it  $O(\log n)$ ).

## Summary of Key Concepts:

- **Hashing:** Converts a key into a hash code, which determines the index in an array.
- **Collision Handling:** Uses either linked lists or tree structures to handle multiple keys that map to the same index.
- **Resizing:** Adjusts the array size as the number of entries grows to maintain efficiency.

Would you like to explore any particular part of this process in more detail?



How does a HashMap handle collisions?  
 Explain techniques like chaining and open addressing.  
 What happens when multiple keys hash to the same index?  
 Discuss how collisions are resolved and their impact on performance.



A `HashMap` handles collisions—situations where multiple keys hash to the same index—through various techniques to ensure the data structure remains efficient. The most commonly used collision-handling strategies are **chaining** and **open addressing**. Here's a detailed explanation of both, along with their impact on performance.

### 1. Chaining (Used in Java's `HashMap`):

- **How It Works:** In chaining, each index in the underlying array (bucket) holds a reference to a data structure, typically a linked list or, in some implementations, a more complex structure like a balanced tree (since Java 8). If multiple keys hash to the same index, they are all stored in this data structure.
- **Process:**
  - When a key-value pair is inserted, the `HashMap` uses the hash code of the key to determine the index of the bucket.

- If the bucket is empty, the entry is stored directly.
  - If the bucket is already occupied by one or more entries (a collision), the new key-value pair is added to the linked list (or tree) at that bucket.
  - When retrieving a value, the list/tree is searched to find the correct key.
  - Java uses `equals()` to compare keys to distinguish between different objects that have the same hash code.
- **Tree Optimization** (Since Java 8):
    - If the number of entries in a bucket's linked list exceeds a threshold (usually 8), the list is converted into a balanced **red-black tree**. This reduces the search time in that bucket from  $O(n)$  to  $O(\log n)$ .
  - **Advantages:**
    - Chaining allows the `HashMap` to handle collisions without probing other array slots.
    - It's flexible because the linked list/tree can grow as needed without resizing the array immediately.
  - **Disadvantages:**
    - Extra memory is required to store the linked list/tree nodes.
    - If too many collisions occur, the time complexity for retrieval or insertion degrades to  $O(n)$ , especially if the bucket contains many elements in a linked list.
  - **Impact on Performance:**
    - In general, for a well-distributed hash function, the number of collisions is minimal, and operations (insert, delete, retrieve) remain close to  $O(1)$ .
    - If many collisions occur, performance degrades as searching through a long list or tree becomes more expensive. However, Java's conversion to a tree helps mitigate this in extreme cases.

## 2. Open Addressing (Not used in Java's `HashMap`, but common in other hash maps like Python's `dict`):

- **How It Works:** In open addressing, all entries are stored directly in the array itself. If a collision occurs, instead of storing multiple entries in a bucket (chaining), the `HashMap` probes for the next available slot in the array to store the new entry.
- **Types of Probing:**
  - **Linear Probing:** If a collision occurs, the `HashMap` checks the next slot (index + 1). If that slot is occupied, it keeps checking subsequent slots until an empty one is found.
  - **Quadratic Probing:** Instead of checking the next slot linearly, the probing sequence uses a quadratic function to find the next available slot. For example, the steps might be index +  $1^2$ , index +  $2^2$ , index +  $3^2$ , and so on.

- **Double Hashing:** In this technique, a secondary hash function is applied to the key when a collision occurs. The secondary hash function provides the increment step to find the next open slot.

- **Process:**

- When inserting, if the computed index is already occupied, probing starts to find an available index.
- For retrieval, the same probing sequence is used to check where the desired key might have been stored.

- **Advantages:**

- Open addressing avoids the use of extra memory for linked lists or trees.
- For low load factors (ratio of stored entries to total capacity), performance can be faster than chaining because all elements are stored directly in the array.

- **Disadvantages:**

- Open addressing performs poorly when the load factor increases. As the array fills up, more probing is required to find open slots, leading to performance degradation.
- If too many collisions occur, the performance for insertion, search, and deletion degrades to  $O(n)$  as the `HashMap` has to probe multiple positions.
- Deletion can be tricky, as removing an entry could disrupt the probing sequence and make some keys unreachable unless special markers (like tombstones) are used.

- **Impact on Performance:**

- For low load factors (typically less than 0.5), open addressing provides  $O(1)$  average case performance for insertion, deletion, and search.
- When the array becomes nearly full, the number of probes needed increases, leading to worse performance, potentially  $O(n)$  in the worst case.

### 3. Handling Multiple Keys That Hash to the Same Index:

- **In Chaining:**

- When multiple keys hash to the same index, they are stored in a list or tree at that index. Each key in the list/tree is checked for equality using `equals()` to ensure that the correct key-value pair is found.
- In a worst-case scenario, if all keys map to the same bucket, the list could contain all the entries, degrading the performance to  $O(n)$  for search or insertion. However, the use of red-black trees after a threshold helps maintain efficient operations, reducing search time to  $O(\log n)$ .

- **In Open Addressing:**

- When multiple keys hash to the same index, the `HashMap` probes for an empty slot. During retrieval, the same probing sequence is followed to locate the key. If too many keys map to the same index, the `HashMap` may need to probe a long sequence of slots, leading to degraded

performance.

#### 4. Resizing and Load Factor:

- Whether using chaining or open addressing, **resizing** is crucial to maintaining performance. In both techniques, once the load factor exceeds a certain threshold (typically 0.75), the `HashMap` creates a larger array (usually twice the size) and **rehashes** all existing entries into the new array.
- Resizing helps reduce collisions and keeps the average time complexity for operations close to  $O(1)$ .

#### 5. Comparison of Chaining and Open Addressing:

Aspect	Chaining	Open Addressing
<b>Collision Resolution</b>	Stores collisions in linked lists or trees	Uses probing to find the next available slot
<b>Memory Usage</b>	Extra memory for linked lists or trees	All data stored in the array
<b>Performance (Low Load)</b>	$O(1)$ for well-distributed data	$O(1)$ for well-distributed data
<b>Performance (High Load)</b>	Degrades to $O(n)$ for long lists; trees improve it	Degrades to $O(n)$ as the number of probes increases
<b>Load Factor Impact</b>	Not significantly affected, as lists can grow	Highly affected, performance degrades as array fills
<b>Resize Complexity</b>	Requires rehashing all elements during resize	Requires rehashing all elements during resize

#### Summary of Collisions and Their Impact:

- **Chaining:** When collisions occur, entries are stored in linked lists or trees, with performance dropping if the load factor becomes high and many collisions happen in the same bucket. The performance can degrade to  $O(n)$  but is optimized with tree-based structures.
- **Open Addressing:** Probes the next available slots when a collision occurs, leading to performance degradation as the array fills up. Operations can degrade to  $O(n)$  in high-collision scenarios, especially at high load factors.

Chaining is generally more flexible and handles higher load factors better, whereas open addressing can be faster for low-load cases but is more sensitive to array capacity limits.



What is the significance of the `hashCode()` and `equals()` methods?  
Discuss how they affect key storage and retrieval.



The `hashCode()` and `equals()` methods are critical in ensuring the correct functioning of `HashMap` (and other hash-based data structures) in Java. Their significance lies in how keys are stored and retrieved. The relationship between these two methods determines whether two objects are treated as equal and ensures that they are placed and found in the correct bucket during storage and retrieval operations. Here's a detailed breakdown of their importance:

## 1. The `hashCode()` Method:

- **Purpose:** The `hashCode()` method provides an integer representation (hash code) for the object, which is used by hash-based collections like `HashMap` to determine the index at which the key-value pair should be stored.

- **How It Works:**

- When an object is inserted into a `HashMap` as a key, its `hashCode()` method is called to compute its hash code.
- The hash code is then used to calculate the index in the underlying array (bucket). Typically, this is done using the formula:

```
java
```

```
index = hashCode() % array.length;
```

- If two keys have the same hash code, they will end up in the same bucket, leading to a collision. This is where the `equals()` method comes into play (explained later).
- **Significance:**
  - A good `hashCode()` function ensures a uniform distribution of keys across the `HashMap`, minimizing the likelihood of collisions and ensuring efficient access.
  - If `hashCode()` produces poor distribution (i.e., too many objects produce the same hash code), it leads to more collisions, which slows down retrieval operations.
- **Contract with `equals()`:**
  - Objects that are considered equal (based on the `equals()` method) **must** have the same hash code.
  - The reverse is **not necessarily true**: two objects with the same hash code **do not** have to be equal. This is essential because collisions can occur.

- **Example:**

```
java
```

```
@Override
public int hashCode() {
    return Objects.hash(id, name);
}
```

In this example, the hash code is computed based on the `id` and `name` fields of an object. Two objects with the same `id` and `name` will have the same hash code.

## 2. The `equals()` Method:

- **Purpose:** The `equals()` method defines the logical equality of two objects. It is used to compare objects for equality when there is a collision (i.e., when two objects hash to the same bucket).
- **How It Works:**
  - When two keys have the same hash code and fall into the same bucket, the `equals()` method is used to compare the keys to determine whether they are logically the same.
  - If `equals()` returns `true`, the `HashMap` considers the two objects as the same key.
  - If `equals()` returns `false`, even though the hash codes match, the keys are considered distinct, and both will be stored in the same bucket (as separate entries in a list or tree).
- **Significance:**
  - The `equals()` method ensures the correctness of the `HashMap` when collisions occur. Without it, the `HashMap` might incorrectly treat two different objects as identical just because they have the same hash code.
  - It allows the `HashMap` to distinguish between different keys that may have the same hash code.
  - The method must adhere to the `equals()` contract: it should be reflexive, symmetric, transitive, consistent, and return `false` for `null` comparisons.
- **Contract with `hashCode()`:**



- If two objects are equal according to the `equals()` method, they **must** have the same hash code. This is critical to maintaining the correctness of the `HashMap`.
- Violating this contract (i.e., two objects being equal but having different hash codes) will lead to unpredictable behavior in the `HashMap`. For example, two equal objects could be stored in different buckets, causing lookups to fail.

- **Example:**

```
java

@Override
public boolean equals(Object obj) {
    if (this == obj) return true;
    if (obj == null || getClass() != obj.getClass()) return false;
    MyObject that = (MyObject) obj;
    return id == that.id && name.equals(that.name);
}
```

In this example, two objects are considered equal if their `id` and `name` fields are equal.

### 3. How They Affect Key Storage and Retrieval:

- **Storage (`put()` operation):**

1. When a key-value pair is inserted using the `put()` method, the `hashCode()` of the key is computed to determine the bucket (index in the array) where the key-value pair should be stored.
2. If the bucket is empty, the key-value pair is placed directly in that bucket.
3. If the bucket already contains one or more entries (i.e., a collision has occurred), the `equals()` method is used to check if any of the existing keys in that bucket are equal to the new key.
  - If an equal key is found, the new value replaces the old value for that key.
  - If no equal key is found, the new key-value pair is added to the bucket (in a list or tree structure).

- **Retrieval (`get()` operation):**

1. When retrieving a value using the `get()` method, the `hashCode()` of the key is computed to locate the bucket where the key-value pair might be stored.
2. If the bucket is empty, the `get()` method returns `null` (the key doesn't exist in the map).
3. If the bucket contains entries, the `equals()` method is used to compare the keys in that bucket with the provided key.

- If an equal key is found, the corresponding value is returned.
- If no equal key is found, `null` is returned.

- **Example of How `hashCode()` and `equals()` Work Together:** Consider two objects, `obj1` and `obj2`, with the same hash code but not equal according to the `equals()` method:

```
java
```

```
MyObject obj1 = new MyObject(1, "Alice");
MyObject obj2 = new MyObject(2, "Bob");
```

- `hashCode()` for `obj1` and `obj2` may both return the same value, say `42`, due to some hash code collision.
- However, `equals(obj1, obj2)` returns `false` because their `id` fields are different.
- In a `HashMap`, both `obj1` and `obj2` will be stored in the same bucket but treated as distinct keys.

#### 4. Potential Issues if `hashCode()` and `equals()` Are Incorrectly Implemented:

- **Inconsistent `hashCode()`:** If two equal objects return different hash codes, they will be placed in different buckets, causing the `HashMap` to fail when retrieving a value using one of those objects as a key.
  - For example, if you put a key into a `HashMap` and later change its state so that its `hashCode()` changes, you may no longer be able to retrieve the value, because the `HashMap` will search in the wrong bucket.
- **Inconsistent `equals()`:** If `equals()` is not implemented properly (e.g., violates symmetry or transitivity), the `HashMap` might store duplicate keys that should have been considered equal, leading to incorrect behavior.

#### 5. Key Points:

- `hashCode()` determines the bucket in which the key will be stored.
- `equals()` is used to check for logical equality between keys, especially when collisions occur.
- If two objects are equal according to `equals()`, they **must** have the same hash code. If this is violated, the `HashMap` will behave unpredictably.
- A well-designed `hashCode()` function ensures good distribution of keys across the buckets, while a correct `equals()` implementation ensures the correct comparison of keys.

These methods are essential for the efficiency and correctness of `HashMap` operations, and improper implementations can lead to bugs, data loss, or degraded performance.



How does a HashMap differ from other data structures like TreeMap or LinkedHashMap?  
Compare and contrast their use cases and performance.



A `HashMap`, `TreeMap`, and `LinkedHashMap` are all implementations of the `Map` interface in Java, but they differ significantly in how they store data, their performance characteristics, and their use cases. Here's a detailed comparison of these data structures:

## 1. HashMap:

- **Underlying Data Structure:** Array with linked lists (or red-black trees in Java 8+) for collision handling.
- **Order of Elements:** `HashMap` **does not maintain any order** of the keys. The order of the elements when iterating is unpredictable and can change if the map is modified.
- **Performance:**
  - **Time Complexity:**
    - Average-case:  $O(1)$  for `put()`, `get()`, and `remove()`.
    - Worst-case:  $O(n)$  in the case of many hash collisions, but this is mitigated in Java 8+ with tree-based buckets (reducing the worst-case to  $O(\log n)$ ).
  - Very fast for most general-purpose operations, especially when order is not a concern.
- **Use Cases:**
  - Suitable when **fast access, insertion, and deletion** are required and **order is irrelevant**.
  - Example use case: A cache system where the order of elements is not important, but fast retrieval is critical.

## 2. TreeMap:

- **Underlying Data Structure:** A **Red-Black Tree** (a type of self-balancing binary search tree).
- **Order of Elements:** `TreeMap` **maintains a sorted order** based on the natural ordering of the keys (via `Comparable`) or a custom order defined by a `Comparator` provided at map construction.
- **Performance:**
  - **Time Complexity:**
    - $O(\log n)$  for `put()`, `get()`, and `remove()` due to the logarithmic height of the red-black tree.
  - Slower than `HashMap` for most operations due to the overhead of maintaining the tree structure.
- **Use Cases:**

- Suitable when **sorted order of keys** is required or useful. For example, if you need to iterate over the map in key order or perform range queries (e.g., finding all keys between a certain range).
- Example use case: A navigation system where locations are stored as keys and need to be sorted based on their names or distances.

### 3. LinkedHashMap:

- **Underlying Data Structure:** Hash table + doubly linked list across the entries.
- **Order of Elements:** `LinkedHashMap` maintains the **insertion order** of the elements. You can also configure it to maintain access order (the order in which elements are accessed).
- **Performance:**
  - **Time Complexity:**
    - $O(1)$  for `put()`, `get()`, and `remove()` on average, the same as `HashMap`.
    - Slightly slower than `HashMap` due to the overhead of maintaining the linked list for ordering.
  - Access order can be beneficial for implementing features like a **least-recently-used (LRU) cache**.
- **Use Cases:**
  - Suitable when you need the **insertion order** or **access order** to be preserved during iteration. It's commonly used when you want the behavior of `HashMap` but with predictable iteration order.
  - Example use case: An LRU cache where recently accessed items are pushed to the end of the list, and the oldest elements are evicted when the cache reaches its size limit.

### 4. Comparison Table:

Feature/Characteristic	HashMap	TreeMap	LinkedHashMap
<b>Underlying Data Structure</b>	Hash table + linked list/tree (for collisions)	Red-Black Tree (self-balancing)	Hash table + doubly linked list
<b>Order of Elements</b>	No order (unpredictable)	Sorted by key (natural or custom order)	Insertion order or access order
<b>Time Complexity (put/get/remove)</b>	$O(1)$ average, $O(n)$ worst (but $O(\log n)$ in Java 8+)	$O(\log n)$	$O(1)$ average
<b>Memory Overhead</b>	Low (array + linked list/tree for collisions)	Higher (tree structure)	Moderate (hash table + linked list)
<b>Best Use Cases</b>	Fast access, insertion, deletion where order is irrelevant	Sorted maps, range queries, or when sorted traversal of keys is needed	Use when predictable iteration order is needed (insertion or access order)
<b>Common Scenarios</b>	General-purpose usage, caching, lookup tables	Sorted collections, ordered maps	Caches, preserving insertion order, implementing LRU behavior

## 5. Performance and Use Cases in More Detail:

### HashMap:

- **Performance:**
  - Best for scenarios where **order does not matter** and operations need to be fast. `HashMap` is highly efficient for lookups and insertions when the key distribution is uniform.
  - **Drawbacks:** You cannot iterate through the entries in a predictable order, so if you need the order of insertion or sorting, `HashMap` is not suitable.
- **Use Case Example:**
  - A mapping of user IDs to user objects where you don't care about the order, but you need fast lookups and updates.

### TreeMap:

- **Performance:**
  - `TreeMap` is slower than `HashMap` for all basic operations because of the overhead of maintaining a balanced tree. However, it is the best choice when you need the **keys to be sorted** at all times.
  - **Drawbacks:** The tree structure adds more overhead than a hash table, so performance for insertions and lookups is slower ( $O(\log n)$ ).
- **Use Case Example:**
  - Storing data in a dictionary where the keys need to be in a **sorted order**, such as storing stock prices by timestamp where you may need to query by a range of timestamps.

### LinkedHashMap:

- **Performance:**
  - Performs similarly to `HashMap`, but with a slight overhead due to the maintenance of the doubly linked list for maintaining order. This overhead is typically minimal in most use cases.
  - **Drawbacks:** Marginally slower than `HashMap` for insertion and deletion due to the linked list maintenance, but it's very useful when insertion or access order matters.
- **Use Case Example:**
  - A **cache system** where you want to maintain the order of elements by how recently they were accessed. You can use it to create an LRU cache, which automatically evicts the least recently accessed item once the cache reaches a certain size.

## 6. Choosing the Right Data Structure:

- **Use** `HashMap` when:
  - You need the fastest performance for key-based lookups, insertions, and deletions.
  - You do not care about the order in which entries are stored or retrieved.
  - Example: A hash-based cache or a mapping of IDs to object references.
- **Use** `TreeMap` when:
  - You need the keys to be in **sorted order**, either natural or based on a custom comparator.
  - You need to perform operations that involve **range queries**, like finding all keys within a certain range.
  - Example: A sorted log of timestamps and corresponding events.
- **Use** `LinkedHashMap` when:
  - You need the elements to maintain **insertion order** or **access order** during iteration.
  - You want to implement a **cache** with predictable eviction behavior, such as an LRU cache.
  - Example: Maintaining a log of recently accessed web pages in the order they were accessed.

## Conclusion:

- `HashMap` is best for fast, unordered access.
- `TreeMap` is ideal when you need a sorted order of keys and can afford the slightly slower performance.
- `LinkedHashMap` is useful when you want predictable iteration order (insertion or access order) with performance similar to `HashMap`.

Choosing between these data structures depends on your specific requirements regarding performance and the order of elements.

ChatGPT can make mistakes. Check important info.