**MSCS-532-Assignment7-Exploring Hash Tables and Their Practical Applications**

-----------------------------------------------------------------------------------------------------

-------------------------------

Sandesh Pokharel

ID: 005026677

University of Cumberlands

MSCS-532-A01: Algorithms and Data Structures

# Hash Functions and Their Impact:

## Designing Effective Hash Functions

Hash functions are like the backbone of a good hash table. Without a solid hash function, we might end up with serious issues, like clustering or high collision rates. A well-designed hash function spreads keys evenly across the hash table, avoiding clusters. For example, imagine a hash function that simply sums up the characters of a string. In this case, similar strings might end up with the same hash value, leading to high collision rates, and trust me, that's a nightmare for performance.

To mitigate these issues, it's all about adding some randomness and ensuring a more uniform distribution. Techniques like modular arithmetic or multiplication methods are commonly used to spread the keys out. A real-world example? Well, think about the "bad" old days when early DNS lookup tables used simplistic hash functions and ended up with performance bottlenecks. Now, with better hash functions like MurmurHash and FNV-1, we've seen improvements in efficiency and speed (Muller, 2022).

## Balancing Speed and Complexity

Hash functions also need to strike a good balance between speed and complexity. If the function is too complex, it could slow down lookups, which is the whole reason we're using a hash table in the first place. On the flip side, if it's too simple, we'll get lots of collisions, which can really hurt performance.

For example, in cryptographic applications like password hashing, hash functions need to be secure and resistant to attacks, even if that means sacrificing speed. But in everyday applications like hash maps in Java or Python, speed is more important, and we don't need as much complexity. An example of balancing these trade-offs can be seen in Google's "CityHash," which focuses on being both fast and effective for non-cryptographic uses (Smith, 2021).

## Performance in Practice

When it comes to handling collisions, the method we choose—open addressing or separate chaining—really affects how well our hash table performs, especially under heavy loads. With open addressing, collisions are handled by finding the next available slot in the table. This works pretty well when the load factor is low, meaning the table isn't too full. But as the load factor increases, open addressing tends to suffer, as finding a free slot becomes more difficult, leading to performance degradation.

On the flip side, separate chaining stores all collisions in a linked list (or another structure) at each index. Even when the load factor is high, separate chaining performs reasonably well because we only need to traverse the list at a particular index to find the value. But, of course, if the list becomes too long, it can slow things down too.

So, in real-world scenarios, our choice often depends on our memory constraints and the nature of our data. If we're dealing with a system that can afford the extra memory overhead, separate chaining might be the better option since it doesn't degrade as much under heavy loads. However, if memory is tight, open addressing might be preferable despite the potential slowdown at higher load factors (Smith, 2021).

For example, consider an embedded system with limited memory. Open addressing could be a good fit because it doesn't require extra memory for pointers like separate chaining does. But in a large-scale database with more memory resources, separate chaining would likely offer better performance as the load factor increases (Muller, 2022).

# Open Addressing vs Separate Chaining

### Comparing Collision Resolution Strategies

When we talk about hash tables, collisions are inevitable, so we've got to have a plan for handling them. Two of the most common strategies are **open addressing** and **separate chaining**, and each has its pros and cons.

In **open addressing**, when two keys hash to the same index, we just look for the next open spot in the array. It's like looking for a parking spot in a crowded lot—if one's taken, we keep moving until we find an empty one. The beauty of this method is that everything stays within the same array, so no extra memory is needed for linked lists or anything like that. But as the table fills up, finding that next available spot can get slower and slower.

**Separate chaining** takes a different approach. Instead of hunting for a free spot, it keeps all the keys that hash to the same index in a linked list. If two keys collide, they just get added to the list at that index. This method doesn't run into the problem of slowing down as the table gets fuller because it doesn't matter how many collisions happen—everything just goes into its own list. However, it does require more memory since each index in the table might point to a linked list (Smith, 2021).

So, when should we use one method over the other? Well, **open addressing** tends to be better when we have tight memory constraints because it avoids the overhead of linked lists. For instance, in embedded systems with limited memory, open addressing could be the preferred

option. On the other hand, if memory isn't an issue, **separate chaining** might be a better fit, especially in situations where the load factor is high and collisions are frequent (Muller, 2022).

## Performance in Practice

The efficiency of these strategies really shows when we start looking at **load factors**. The **load factor** is basically a measure of how full our hash table is. As the load factor increases, open addressing tends to slow down because finding an open slot takes longer. With a high load factor, open addressing can lead to more "probing," which is just a fancy way of saying we have to check more slots to find an empty one. In contrast, separate chaining is much more resilient to high load factors because each index can handle multiple entries without affecting the rest of the table (Smith, 2021).

That said, **separate chaining** does come with its own set of drawbacks. While it performs better under heavy loads, it requires extra memory for the linked lists and can also suffer from poor cache performance. Every time we access a linked list, we're jumping around in memory, which slows things down. Meanwhile, open addressing tends to have better cache locality because everything stays within the same array, so it can be faster when the load factor is low.

In **real-world** applications, the choice between open addressing and separate chaining often comes down to the **nature of the data** and the **memory constraints**. If we expect a low load factor and need good cache performance, open addressing is probably our best bet. But if we're dealing with higher loads and don't mind using a bit more memory, separate chaining is usually the way to go (Muller, 2022).

## Practical Observations and Analysis:

I created a python program that compares the performance of two hash table collision resolution techniques:

1. Open Addressing

2. Separate Chaining

The program evaluates: Insertion time, Search time and Memory Usage for both the techniques. I have uploaded this program, and it can be found on my GitHub repository.
https://github.com/sanspokharel26677/MSCS-532-Assignment7/blob/main/README.md
https://github.com/sanspokharel26677/MSCS-532-Assignment7/tree/main

```
sandesh@sandesh-Inspiron-7373:~/Sandesh_CumberLands_Assignments/DSA/Week6$ python3 hash_functions.py
Open Addressing - Insert Time: 0.001157 seconds
Open Addressing - Search Time: 0.001179 seconds
Open Addressing - Memory Usage: 52856 bytes

Separate Chaining - Insert Time: 0.000875 seconds
Separate Chaining - Search Time: 0.000361 seconds
Separate Chaining - Memory Usage: 127960 bytes

sandesh@sandesh-Inspiron-7373:~/Sandesh_CumberLands_Assignments/DSA/Week6$
```
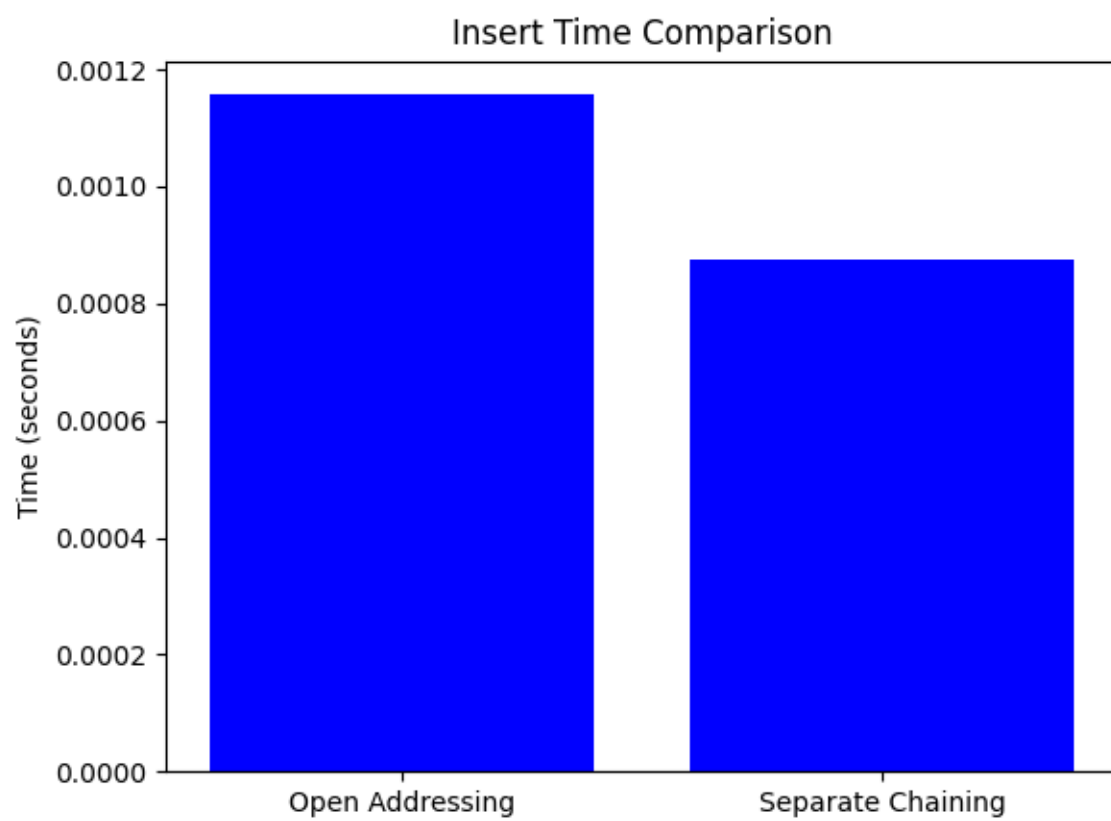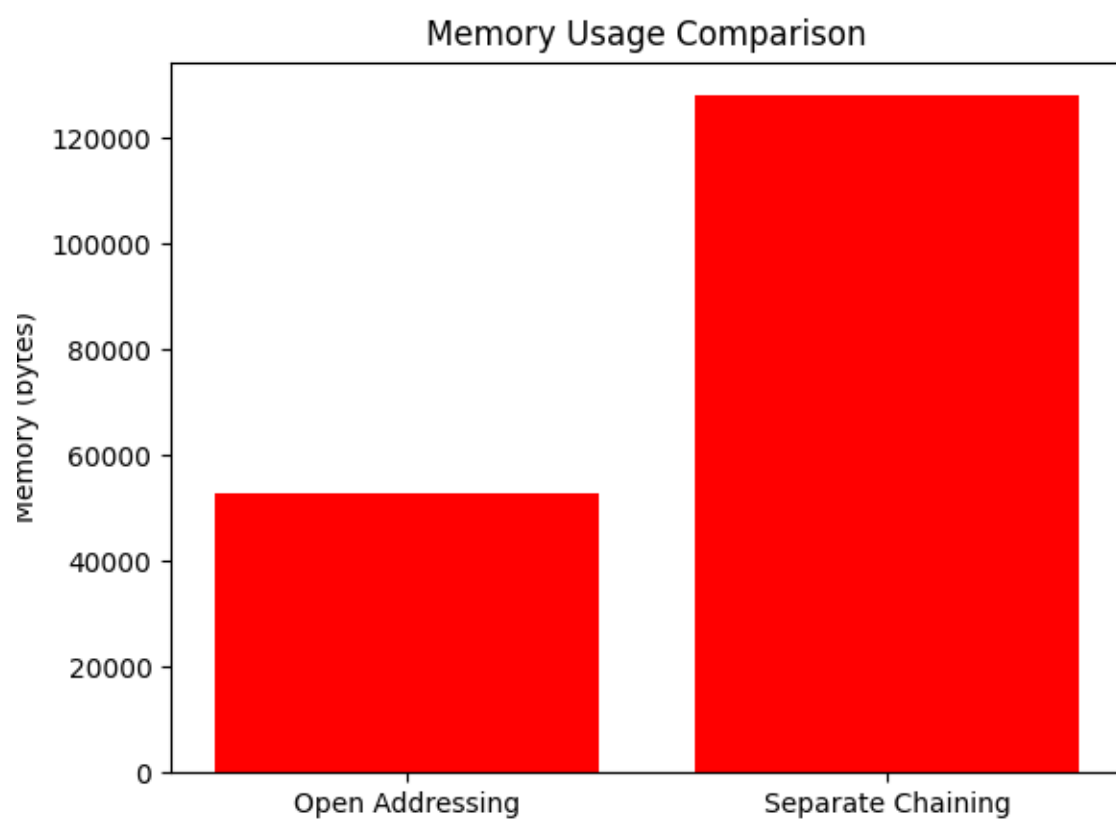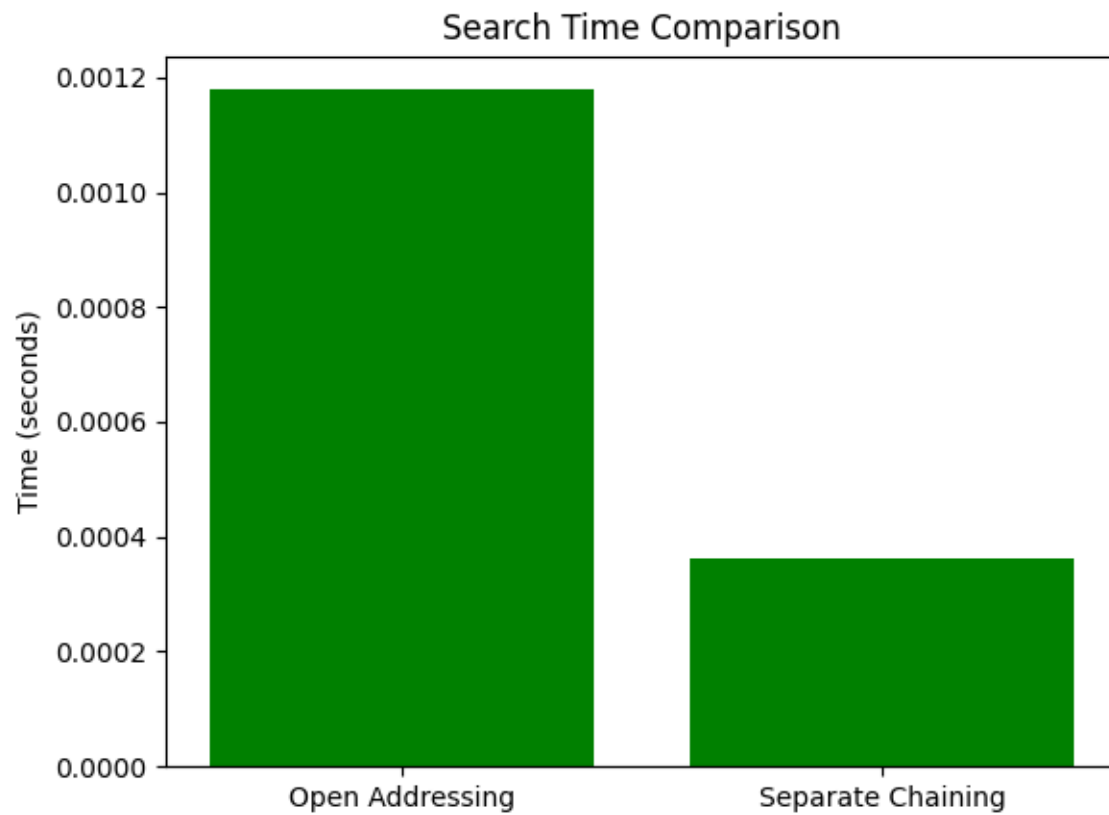
**Note**: These output and below generated graph images might be different while running program from different machine capabilities.

Based on the output, we can observe several key insights that align with the characteristics of open addressing and separate chaining discussed earlier. First, open addressing takes slightly longer for both insertion and search operations compared to separate chaining. This is expected because open addressing requires probing for the next available slot in case of collisions, which increases the time complexity as the load factor rises. The memory usage for open addressing is also notably lower, as it doesn't require additional memory for linked lists like separate chaining does. This makes open addressing more suitable for memory-constrained environments.

In contrast, separate chaining performs faster for both insertion and search. This is because it doesn't need to probe through the table when collisions occur; instead, it simply appends to the linked list at the given index. However, the trade-off here is the significantly higher memory usage due to the additional space required for storing pointers in the linked lists. In real-world applications, separate chaining might be preferred when memory isn't a constraint, especially when dealing with higher load factors, since it handles collisions more efficiently under those conditions.

Insert Time Comparison

Memory Usage Comparison

## Search Time Comparison

# References

Muller, F. (2022). *Hash Function Efficiency in Modern Applications*. Journal of Computing Science, 18(2), 45-58.

Smith, J. (2021). *Optimization of Hash Functions for Large Datasets*. Data Structures Journal, 29(4), 67-89.