

Hash Tables in C++: A Refresher for C Programmers with AS/CR Style

Your Name

Overview and Motivation

A **hash table** is a data structure that lets you store and retrieve key-value pairs in nearly constant time on average. If you've ever used a dictionary or an associative array, you've used something like a hash table.

Why Are Hash Tables Important?

- Fast lookup: $O(1)$ average case for insertion, deletion, and search.
- Flexible: Can store almost anything if you define how to hash it.
- Foundation for many abstractions: caches, symbol tables, compiler internals, network routing tables, etc.
- Crucial for performance-oriented design—particularly important when building data-centric systems or modeling abstract relationships efficiently (e.g., graph traversal, sparse matrices).

Intuitive Analogy (CR Friendly)

Imagine you're in a library where books are stored not by order, but by hashing the title into a number that tells you exactly which locker to open. If the locker is empty, place the book there. If there's a collision (another book already in the locker), you have to handle it by some clever strategy.

Key Insight: Hash Function + Array = Superpower

Hash tables are made of two simple parts:

1. A **hash function** $h(key)$: maps the key to an array index.
2. An **array** (a contiguous block of memory): stores key-value pairs at these indices.

The power of a hash table lies in how cleverly you handle collisions and design your hash function.

Struct Refresher (AS Foundation)

In C, a **struct** is a way to group variables (of possibly different types) under one name.

```
1 struct Book {  
2     int id;  
3     const char* title;  
4     float price;  
5 };
```

Listing 1: Basic struct definition in C-style

In C++, this becomes richer with constructors and methods.

```

1 struct Book {
2     int id;
3     std::string title;
4     float price;
5
6     // Constructor
7     Book(int i, std::string t, float p) : id(i), title(t), price(p) {}
8 };

```

Listing 2: C++ struct with constructor

You'll see structs used in hash table nodes (e.g., key-value pairs).

Pointers Refresher (CR-Style Insight)

Pointers let you reference memory locations directly, which is necessary when managing dynamic memory (especially with collision chains or rehashing).

```

1 int a = 10;
2 int* ptr = &a;
3 std::cout << *ptr; // Outputs 10

```

Listing 3: Simple pointer usage

Important: In hash tables, we often use pointers to link items together (like in linked lists), or to allocate storage dynamically for flexible growth.

Hash Table Components (AS Map)

- **Key:** What you're indexing by (e.g., name, number).
- **Value:** What you're storing.
- **Bucket:** The array slot where items get stored.
- **Hash function:** Converts the key to an index.
- **Collision resolution strategy:** What to do if two keys hash to the same index.

Common Collision Resolution Techniques

- **Chaining:** Each bucket holds a linked list (or dynamic array) of entries.
- **Open addressing:** Look for the next available slot (e.g., linear probing, quadratic probing, double hashing).

Coming Up Next

In the next module, we'll build a minimal working hash table in C++ using:

- Arrays of pointers to linked lists (for chaining).
- A simple hash function for strings or integers.
- Custom structs for key-value pairs.

Preview: You'll build something like this:

```

1 // HashTable
2 // |--- Bucket 0: NULL
3 // |--- Bucket 1: [ ("Alice", 100) -> ("Bob", 95) ]
4 // |--- Bucket 2: NULL
5 // |--- ...

```

Design Overview

We will build a hash table that:

- Stores `std::string` keys and `int` values.
- Uses a basic hash function.
- Uses **chaining** with linked lists to handle collisions.

High-Level Structure

- A fixed-size array of `Node*`, where each bucket is a pointer to the head of a linked list.
- Each `Node` holds a key, value, and pointer to the next node.

Step 1: Define a Node Struct

```

1 struct Node {
2     std::string key;
3     int value;
4     Node* next;
5
6     Node(std::string k, int v) : key(k), value(v), next(nullptr) {}
7 };

```

Listing 4: Struct for Linked List Node

AS insight: We're grouping key-value pairs with a next pointer to form a singly linked list.

CR insight: This structure will allow us to "chain" multiple key-value pairs at the same array index if collisions occur.

Step 2: Define the Hash Table Class

```

1 class HashTable {
2 private:
3     static const int TABLE_SIZE = 10;
4     Node* table[TABLE_SIZE];
5
6     int hashFunction(std::string key);
7
8 public:
9     HashTable();
10    ~HashTable();
11
12    void insert(std::string key, int value);
13    int search(std::string key);
14    void remove(std::string key);
15 };

```

Listing 5: Basic Hash Table Class Skeleton

Step 3: Constructor and Destructor

```
1 HashTable::HashTable() {
2     for (int i = 0; i < TABLE_SIZE; ++i) {
3         table[i] = nullptr;
4     }
5 }
6
7 HashTable::~~HashTable() {
8     for (int i = 0; i < TABLE_SIZE; ++i) {
9         Node* entry = table[i];
10        while (entry != nullptr) {
11            Node* prev = entry;
12            entry = entry->next;
13            delete prev;
14        }
15    }
16 }
```

Why do this?

- We initialize all buckets to `nullptr`.
- In the destructor, we walk through each bucket and free any nodes. This avoids memory leaks.

Step 4: A Simple Hash Function

```
1 int HashTable::hashFunction(std::string key) {
2     int hash = 0;
3     for (char c : key) {
4         hash += c;
5     }
6     return hash % TABLE_SIZE;
7 }
```

Listing 6: Hash Function for Strings

Note: This is a weak hash, but it works for small demos.

Step 5: Insert Operation

```
1 void HashTable::insert(std::string key, int value) {
2     int index = hashFunction(key);
3     Node* newNode = new Node(key, value);
4     newNode->next = table[index];
5     table[index] = newNode;
6 }
```

This adds the new node to the front of the list at `table[index]`.

Step 6: Search Operation

```
1 int HashTable::search(std::string key) {
2     int index = hashFunction(key);
3     Node* entry = table[index];
4 }
```

```

5     while (entry != nullptr) {
6         if (entry->key == key) return entry->value;
7         entry = entry->next;
8     }
9     return -1; // Key not found
10 }

```

Step 7: Remove Operation

```

1 void HashTable::remove(std::string key) {
2     int index = hashFunction(key);
3     Node* entry = table[index];
4     Node* prev = nullptr;
5
6     while (entry != nullptr && entry->key != key) {
7         prev = entry;
8         entry = entry->next;
9     }
10
11     if (entry == nullptr) return; // Not found
12
13     if (prev == nullptr) {
14         table[index] = entry->next; // Remove head
15     } else {
16         prev->next = entry->next;
17     }
18
19     delete entry;
20 }

```

Test It All Together

```

1 int main() {
2     HashTable ht;
3     ht.insert("Alice", 100);
4     ht.insert("Bob", 95);
5     ht.insert("Charlie", 85);
6
7     std::cout << "Alice:_" << ht.search("Alice") << std::endl;
8     std::cout << "Bob:_" << ht.search("Bob") << std::endl;
9
10    ht.remove("Bob");
11    std::cout << "Bob_after_removal:_" << ht.search("Bob") << std::endl;
12
13    return 0;
14 }

```

Listing 7: Main function to test the hash table

CR Wrap-up: Hidden Powers

By understanding this structure, you now know how to:

- Use linked data structures with dynamic memory.

- Control memory layout and access.
- Prepare for advanced modeling like sparse data graphs, lookup acceleration, or key-indexed maps.

Better Hash Functions

In Module 2, our hash function was:

```
1 int hash = 0;
2 for (char c : key) {
3     hash += c;
4 }
5 return hash % TABLE_SIZE;
```

Why it's bad (AS/CR dual view):

- It treats permutations of characters the same ("abc" == "cab").
- It increases collision probability for short ASCII strings.
- It doesn't scale well when 'TABLE_SIZE' increases.

Improved Version: DJB2 Hash

The DJB2 hash function is simple, fast, and surprisingly effective.

```
1 unsigned long HashTable::hashFunction(std::string key) {
2     unsigned long hash = 5381;
3     for (char c : key) {
4         hash = ((hash << 5) + hash) + c; // hash * 33 + c
5     }
6     return hash % TABLE_SIZE;
7 }
```

This gives better key distribution and sets you up for future hashing hardware-accelerated logic.

Dynamic Resizing (a.k.a. Rehashing)

Problem: A fixed table size means eventual performance decay (more collisions, longer chains).

Solution: Track the load factor $\alpha = \frac{\text{num_entries}}{\text{table_size}}$. When $\alpha > 0.75$, double the size and rehash.

Steps to Resize

1. Store the old table.
2. Allocate a new table (2x size).
3. Reinsert every key-value pair into the new table.
4. Free the old table.

We'll cover the implementation in a later advanced module to stay focused on concepts.

Why `std::vector` is Your Friend

Instead of a fixed array:

```
1 Node* table[TABLE_SIZE];
```

We use:

```
1 std::vector<Node*> table;
```

Advantages:

- Auto-resizing (with `resize()`).
- Cleaner memory handling.
- Essential for GPU-style data chunking (next section).

Update Constructor with Vector

```
1 HashTable::HashTable() {  
2     table.resize(TABLE_SIZE, nullptr);  
3 }
```

Spatial Memory Intuition (CR Spatial Insight)

Your hash table's bucket array is like a 1D grid in RAM. With **chaining**, each cell points to a mini list. In GPU terms:

- The table is **device global memory**.
- Each thread can operate on a disjoint bucket or chain.
- All operations must respect memory access rules.

CUDA Insight: Buckets can be parallelized. Each block/thread gets:

- A bucket index.
- A pointer to walk through the chain.

This maps cleanly to warp-style access patterns.

Designing with Parallelism in Mind

Even though we're still in CPU-C++, it pays to design the interface so that:

- Insertions/searches can occur in parallel buckets.
- Hash functions are pure (no state/memory dependencies).
- Data locality is improved (you'll cache better).

Your New Optimized Table Class (Sketch)

```
1 class HashTable {
2 private:
3     std::vector<Node*> table;
4     int size;
5
6     unsigned long hashFunction(std::string key);
7
8 public:
9     HashTable(int tableSize = 10);
10    void insert(std::string key, int value);
11    int search(std::string key);
12    void remove(std::string key);
13 };
```

AS Parallel-Ready Abstraction:

- Input: key
- Compute: index
- Operate: on table[index] independently

Coming Up in Module 4

You'll begin designing:

- A rehashable, vectorized hash table.
- Hooks for benchmarking insert/search speed.
- A "simulated GPU kernel" for parallel insertions (emulated with threads or loops).

Bonus: We'll preview how to port the core structure to a CUDA kernel-friendly format (e.g., SoA/CoA, global vs. shared memory constraints).

Dynamic Resizing (a.k.a. Rehashing)

Your table should grow when too many items crowd into a small number of buckets (i.e., high load factor).

When to Rehash?

Let:

$$\alpha = \frac{\text{num.elements}}{\text{table.size}}$$

Trigger rehash when $\alpha > 0.75$.

Rehashing Steps

1. Save the old table.
2. Create a new table with size = $2 \times$ old size.
3. Reinsert each element.
4. Delete old nodes to avoid memory leaks.

Implementation Details

New member variables:

```
1 int capacity;          // total slots
2 int count;             // total elements
```

Add this to ‘insert()’:

```
1 if ((float)count / capacity > 0.75) {
2     rehash();
3 }
```

Rehash implementation:

```
1 void HashTable::rehash() {
2     int old_capacity = capacity;
3     std::vector<Node*> old_table = table;
4
5     capacity *= 2;
6     table.clear();
7     table.resize(capacity, nullptr);
8     count = 0;
9
10    for (int i = 0; i < old_capacity; ++i) {
11        Node* entry = old_table[i];
12        while (entry) {
13            insert(entry->key, entry->value); // reinsert into new table
14            Node* prev = entry;
15            entry = entry->next;
16            delete prev;
17        }
18    }
19 }
```

Note: We avoid directly copying pointers—this preserves hashing consistency with the new size.

Performance Measurement Hooks

Use the standard C++ ‘`chrono`’ library:

```
1 #include <chrono>
2
3 auto start = std::chrono::high_resolution_clock::now();
4
5 // Perform insertions/searches
6
7 auto end = std::chrono::high_resolution_clock::now();
8 std::chrono::duration<double> duration = end - start;
9 std::cout << "Time:␣" << duration.count() << "s\n";
```

This gives you fine-grained timing for benchmarking insert/search.

Simulating Parallelism (Pre-CUDA Mental Model)

Imagine multiple threads inserting into the hash table. What happens?

Challenge: Two threads inserting into the same bucket may cause race conditions.

CR Insight: Lockless Chaining Design

A CUDA-friendly design requires:

- Preallocated node storage (pool).
- Atomic inserts (e.g., lock-free append).
- Per-bucket synchronization (e.g., fine-grained locks or atomics).

CPU Emulation Idea:

Use ‘std::thread’ to simulate parallel insertions:

```
1 void threadedInsert(HashTable* ht, std::string key, int val) {
2     ht->insert(key, val);
3 }
4
5 std::thread t1(threadedInsert, &ht, "Alice", 100);
6 std::thread t2(threadedInsert, &ht, "Bob", 200);
7
8 t1.join();
9 t2.join();
```

BUT: You must add locks around ‘insert()’ if you do this!

Preparing for GPU-Style Design (Preview)

On a GPU:

- Buckets must be stored in device memory.
- Each thread gets a bucket index (e.g., via hash).
- You’ll need:
 - Flat memory representation (no pointers).
 - Atomic memory ops (e.g., `atomicCAS`).

CR Anticipation: SoA vs. AoS

Instead of:

```
1 struct Node { string key; int val; Node* next; };
```

Think:

- `keys[], values[], next_index[]` → Structure of Arrays (SoA)

This lets you parallelize better in CUDA (coalesced memory access).

Coming in Module 5:

We’ll finalize the full optimized C++ version with:

- Rehashable vectorized table
- Benchmark CLI interface
- Final project file layout
- Then start CUDA-friendly architectural planning!

Final Hash Table Class: Overview

- Fully dynamic and rehashable
- Uses `std::vector` for bucket array
- Benchmarked via `std::chrono`
- Structurally sound for CUDA-style refactoring

Header File: `hash_table.hpp`

```
1  #pragma once
2  #include <vector>
3  #include <string>
4
5  struct Node {
6      std::string key;
7      int value;
8      Node* next;
9      Node(std::string k, int v) : key(k), value(v), next(nullptr) {}
10 };
11
12 class HashTable {
13 private:
14     std::vector<Node*> table;
15     int capacity;
16     int count;
17
18     unsigned long hashFunction(std::string key);
19     void rehash();
20
21 public:
22     HashTable(int initialSize = 10);
23     ~HashTable();
24
25     void insert(std::string key, int value);
26     int search(std::string key);
27     void remove(std::string key);
28 };
```

Implementation File: `hash_table.cpp`

```
1  #include "hash_table.hpp"
2  #include <iostream>
3
4  HashTable::HashTable(int initialSize) : capacity(initialSize), count(0) {
5      table.resize(capacity, nullptr);
6  }
7
8  HashTable::~HashTable() {
9      for (int i = 0; i < capacity; ++i) {
10         Node* entry = table[i];
11         while (entry) {
12             Node* prev = entry;
13             entry = entry->next;
```

```

14         delete prev;
15     }
16 }
17 }
18
19 unsigned long HashTable::hashFunction(std::string key) {
20     unsigned long hash = 5381;
21     for (char c : key)
22         hash = ((hash << 5) + hash) + c;
23     return hash % capacity;
24 }
25
26 void HashTable::rehash() {
27     int old_capacity = capacity;
28     std::vector<Node*> old_table = table;
29
30     capacity *= 2;
31     table.clear();
32     table.resize(capacity, nullptr);
33     count = 0;
34
35     for (int i = 0; i < old_capacity; ++i) {
36         Node* entry = old_table[i];
37         while (entry) {
38             insert(entry->key, entry->value);
39             Node* prev = entry;
40             entry = entry->next;
41             delete prev;
42         }
43     }
44 }
45
46 void HashTable::insert(std::string key, int value) {
47     if ((float)count / capacity > 0.75)
48         rehash();
49
50     int index = hashFunction(key);
51     Node* newNode = new Node(key, value);
52     newNode->next = table[index];
53     table[index] = newNode;
54     ++count;
55 }
56
57 int HashTable::search(std::string key) {
58     int index = hashFunction(key);
59     Node* entry = table[index];
60     while (entry) {
61         if (entry->key == key)
62             return entry->value;
63         entry = entry->next;
64     }
65     return -1;
66 }
67
68 void HashTable::remove(std::string key) {
69     int index = hashFunction(key);
70     Node* entry = table[index];
71     Node* prev = nullptr;
72

```

```

73     while (entry && entry->key != key) {
74         prev = entry;
75         entry = entry->next;
76     }
77
78     if (!entry) return;
79
80     if (!prev) table[index] = entry->next;
81     else prev->next = entry->next;
82
83     delete entry;
84     --count;
85 }

```

Benchmark File: main.cpp

```

1  #include "hash_table.hpp"
2  #include <iostream>
3  #include <chrono>
4
5  int main() {
6      HashTable ht;
7      const int N = 10000;
8
9      auto start = std::chrono::high_resolution_clock::now();
10     for (int i = 0; i < N; ++i) {
11         ht.insert("Key" + std::to_string(i), i);
12     }
13     auto end = std::chrono::high_resolution_clock::now();
14
15     std::chrono::duration<double> diff = end - start;
16     std::cout << "Inserted_" << N << "_items_in_" << diff.count() << "_seconds.\n"
17         ;
18
19     std::cout << "Search_result_for_Key1234:_ " << ht.search("Key1234") << "\n";
20 }

```

Project Structure

```

hash_table_project/
|--- hash_table.hpp
|--- hash_table.cpp
|--- main.cpp
|__ Makefile (optional)

```

Makefile (Optional)

```

1  all:
2      g++ -std=c++17 -O2 main.cpp hash_table.cpp -o hash_table

```

Next Steps (Toward CUDA)

Now that your C++ base is solid, you're ready to:

- Flatten your data layout (SoA instead of pointers)
- Preallocate all memory
- Implement per-bucket insertions in threads
- Use atomic primitives or warp-serial execution

Bonus Thought: Sparse Graphs and Algebraic Models

Your hash table could now easily:

- Represent sparse graphs (edges as keys, weights as values)
- Store fast-access memoization results for recurrence relations
- Index elements by prime-factor "key" patterns

This opens the door to using hash maps in symbolic computation, parallel dataflow models, or even algebraic DSLs you design yourself.