

CS213 2024 Tutorial Solutions

CS213/293 UG TAs

2024

Contents

Tutorial 1	2
Tutorial 2	7
Tutorial 3	12

Tutorial 1

1. The following is the code for insertion sort. Compute the exact worst-case running time of the code in terms of n and the cost of doing various machine operations.

Algorithm 1: Insertion Sort Algorithm

Data: Array A of length n

```
1 for  $j \leftarrow 1$  to  $n - 1$  do
2    $key \leftarrow A[j]$ ;
3    $i \leftarrow j - 1$ ;
4   while  $i \geq 0$  do
5     if  $A[i] > key$  then
6       |  $A[i + 1] \leftarrow A[i]$ ;
7     end
8     else
9       | break;
10    end
11     $i \leftarrow i - 1$ ;
12  end
13   $A[i + 1] \leftarrow key$ ;
14 end
```

Solution: To compute the worst-case running time, we must decide the code flow leading to the worst case. When the if-condition evaluates to false, the control breaks out of the inner loop. The worst case would have the if-condition never evaluates false. This will happen when the input is in a strictly decreasing order.

Counting the operations for the outer iterator j (which goes from 1 to $n - 1$). Consider for each outer loop iteration:

1. Consider the outer loop without the inner loop. Then it has 1 Comparison (Condition in for loop), 1 Increment (1 Assignment + 1 Arithmetic) (Updation in for loop), 3 Assignments (Line 2, 3, 11) and 2 Memory Accesses ($A[j]$ and $A[i+1]$) and 1 Jump. Hence $(C + 2Ar + 4As + 2M + J)$
2. while loop runs for j times (as i goes from $j - 1$ to 0). Without the if-else, the while loop still has to do 1 Comparison (while condition), 1 Decrement (1 Assignment + 1 Arithmetic) ($i \leftarrow i - 1$), 1 Jump, every iteration. Hence $(C + As + Ar + J) * j$
3. The if condition is checked in all iterations of the while loop. It involves 1 Comparison (if) and 1 Memory Access ($A[i]$ in the if condition). if block is executed for all iterations of while in the worst case. if block has 2 Memory Accesses and an Assignment and then a Jump happens to outside the if-else statement. Hence $(C + 3M + As + J) * j$

The total cost would then be (terms are for outer loop, inner loop, if and else):

$$\begin{aligned} & \sum_{j=1}^{n-1} ((C + 2Ar + 4As + 2M + J) + (C + As + Ar + J) * j + (C + 3M + As + J) * j) \\ &= (C + 2Ar + 4As + 2M + J) * (n - 1) + (C + As + Ar + J) * (n * (n - 1) / 2) \\ & \quad + (C + 3M + As + J) * (n * (n - 1) / 2) \end{aligned}$$

This means Insertion Sort is $\mathcal{O}(n^2)$.

(C is Comparison, As is Assignment, Ar is Arithmetic, M is Memory Access, J is Jump)

Note: Some intermediate operations might have been omitted but they do not affect the overall asymptotic performance of the algorithm. The time taken for comparisons, jumps, arithmetic operations, and memory accesses are assumed to be constants for a given machine and architecture.

2. What is the time complexity of binary addition and multiplication? How much time does it take to do unary addition?

Solution: Note that time complexity is measured as a function of the input size since we aim to determine how the time the algorithm takes scales with the input size.

Binary Addition

Assume two numbers A and B. In binary notation, their lengths (number of bits) are m and n. Then the time complexity of binary addition would be $\mathcal{O}(m + n)$. This is because we can start from the right end and add (keeping carry in mind) from right to left. Each bit requires an $\mathcal{O}(1)$ computation since there are only 8 combinations (2 each for bit 1, bit 2, and carry). Since the length of a number N in bits is $\log N$, the time complexity is $\mathcal{O}(\log A + \log B) = \mathcal{O}(\log(AB)) = \mathcal{O}(m + n)$.

Binary Multiplication

Similar to above, but here the difference would be that each bit of the larger number (assumed to be A without loss of generality) would need to be multiplied by the smaller number, and the result would need to be added. Each bit of the larger number would take $\mathcal{O}(n)$ computations, and then m such numbers would be added. So the time complexity would be $\mathcal{O}(m \times \mathcal{O}(n)) = \mathcal{O}(mn)^1$. Following the definition above, the time complexity is $\mathcal{O}(\log A \times \log B) = \mathcal{O}(mn)$.

Side note: How do we know if this is the most efficient algorithm? Turns out, this is NOT the most efficient algorithm. The most efficient algorithm (in terms of asymptotic time complexity) has a time complexity of $\mathcal{O}(n \log n)$ where n is the maximum number of bits in the 2 numbers.

Unary Addition

The unary addition of A and B is just their concatenation. This means the result would have A + B number of 1's. Iterating over the numbers linearly would give a time complexity of $\mathcal{O}(A + B)$ which is linear in input size.

3. Given $f(n) = a_0n^0 + \dots + a_dn^d$ and $g(n) = b_0n^0 + \dots + b_en^e$ with $d > e$, then show that $f(n) \notin \mathcal{O}(g(n))$

Solution: Let us begin by assuming the proposition is False, ergo, $f(n) \in \mathcal{O}(g(n))$. By definition, then, there exists a constant c such that there exists another constant n_0 such that

$$\forall n \geq n_0, f(n) \leq cg(n)$$

. Hence, we have

$$\forall n \geq n_0, a_0n^0 + \dots + a_dn^d \leq cb_0n^0 + \dots + b_en^e$$

$$\forall n \geq n_0, \sum_{i=0}^e (a_i - cb_i)n^i + a_{i+1}n^{i+1} + \dots + a_dn^d \leq 0$$

By definition of limit

$$\lim_{n \rightarrow \infty} \sum_{i=0}^e (a_i - cb_i)n^i + a_{i+1}n^{i+1} + \dots + a_d n^d \leq 0 \\ \implies a_d \leq 0$$

Assuming $a_d > 0$ (since we are dealing with functions mapping from \mathbb{N} to \mathbb{N}), this results in a contradiction; thus, our original proposition is proved.

4. What is the difference between “at” and “[.]” accesses in C++ maps?

Solution: Both accesses will first search for the given key in the map but will behave differently when the key is not present. The at method will throw an exception if the key is not found in the map, but the [] operator will insert a new element with the key and a default-initialized value for the mapped type and will return that default value.

Look at the following illustration:-

```
G++ a.cpp
1 #include <iostream>
2 #include <map>
3 using namespace std;
4
5 void using_square_braces() {
6     map<int, int> temp;
7     temp[0] = 1;
8     cout << temp[0] << endl;
9     cout << temp[1] << endl;
10 }
11
12 void using_at() {
13     map<int, int> temp;
14     temp[0] = 2;
15     cout << temp.at(0) << endl;
16     cout << temp.at(1) << endl;
17 }
18
19 int main() {
20     int choice;
21     cout << "Enter Choice: ";
22     cin >> choice;
23     if(choice == 0) {
24         using_square_braces();
25     }
26     else {
27         using_at();
28     }
29 }
30
```

```
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ a.cpp
kavyagupta@Kavyas-MacBook-Pro Tutorials % ./a.out
Enter Choice: 0
1
0
kavyagupta@Kavyas-MacBook-Pro Tutorials % ./a.out
Enter Choice: 1
2
terminate called after throwing an instance of 'std::out_of_range'
what(): map::at
zsh: abort ./a.out
kavyagupta@Kavyas-MacBook-Pro Tutorials %
```

5. C++ does not provide active memory management. However, smart pointers in C++ allow us the capability of a garbage collector. The smart pointer classes in C++ are

- auto_ptr
- unique_ptr
- shared_ptr
- weak_ptr

Write programs that illustrate the differences among the above smart pointers.

Solution: Memory allocated in heap (using new or malloc) if not de-allocated can lead to memory leaks. Smart pointers in C++ deal with this issue. Broadly, they are classes that store reference counts to the memory they point. When a smart pointer is created, reference count to that memory is increased by one. Whenever a smart pointer (pointing to the same memory) goes out of scope, its destructor is automatically executed, which reduces the reference count of that memory by one and when it hits zero, that memory is deallocated.

- `shared_ptr` allows multiple references to a memory location (or an object)
- `weak_ptr` allows one to refer to an object without having the reference counted
- `unique_ptr` is like `shared_ptr` but it does not allow a programmer to have two references to a memory location
- `auto_ptr` is just the deprecated version of `unique_ptr`

The use for `weak_ptr` is to avoid the circular dependency created when two or more object pointing to each other using `shared_ptr`.

You can see the reference count of a `shared_ptr` using its `use_count()` function.

```

C: shared_ptr.cpp X
1 #include <iostream>
2 #include <memory>
3 using namespace std;
4
5 int main() {
6     shared_ptr<int> P1(new int(4));
7     cout << P1.use_count() << endl;
8     shared_ptr<int> P2 = P1;
9     cout << P2.use_count() << endl;
10    cout << *P2 << endl;
11 }
12
Tutorials -- zsh -- 80x24
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ shared_ptr.cpp
kavyagupta@Kavyas-MacBook-Pro Tutorials % ./a.out
1
2
4
kavyagupta@Kavyas-MacBook-Pro Tutorials %

```

```

C: unique_ptr.cpp X
1 #include <iostream>
2 #include <memory>
3 using namespace std;
4
5 int main() {
6     unique_ptr<int> P1(new int(4));
7     unique_ptr<int> P2 = P1;
8 }
9
Tutorials -- zsh -- 80x24
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ unique_ptr.cpp
unique_ptr.cpp:7:26: error: use of deleted function 'std::unique_ptr<Tp, _Dp>::unique_ptr(const std::unique_ptr<Tp, _Dp>&) [with Tp = int; _Dp = std::default_delete<int>]'
7 |     unique_ptr<int> P2 = P1;
  |     ~~~~~^~~~~
In file included from /opt/homebrew/Cellar/gcc/14.1.0_1/include/c++/14/memory:7
      from unique_ptr.cpp:2:
/opt/homebrew/Cellar/gcc/14.1.0_1/include/c++/14/bits/unique_ptr.h:516:7: note:
declared here
516 |     unique_ptr(const unique_ptr&) = delete;
    |     ~~~~~^~~~~
unique_ptr.cpp:7:26: note: use '-fdiagnostics-all-candidates' to display consid
red candidates
7 |     unique_ptr<int> P2 = P1;
  |     ~~~~~^~~~~
kavyagupta@Kavyas-MacBook-Pro Tutorials %

```

```

C: weak_ptr.cpp X
1 #include <iostream>
2 #include <memory>
3 using namespace std;
4
5 int main() {
6     shared_ptr<int> P1(new int(4));
7     cout << P1.use_count() << endl;
8     weak_ptr<int> P2 = P1;
9     cout << P1.use_count() << endl;
10 }
11
Tutorials -- zsh -- 80x24
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ weak_ptr.cpp
kavyagupta@Kavyas-MacBook-Pro Tutorials % ./a.out
1
1
kavyagupta@Kavyas-MacBook-Pro Tutorials %

```

```

C: cycle_error.cpp X
1 #include <iostream>
2 #include <memory>
3 using namespace std;
4
5 class B;
6
7 class A
8 {
9 public:
10     shared_ptr<B> sP1;
11     A() { cout << "A()" << endl; }
12     ~A() { cout << "~A()" << endl; }
13 };
14
15 class B
16 {
17 public:
18     shared_ptr<A> sP1;
19     B() { cout << "B()" << endl; }
20     ~B() { cout << "~B()" << endl; }
21 };
22
23 int main()
24 {
25     shared_ptr<A> aPtr(new A);
26     shared_ptr<B> bPtr(new B);
27
28     aPtr->sP1 = bPtr;
29     bPtr->sP1 = aPtr;
30
31     cout << aPtr.use_count() << endl;
32     cout << bPtr.use_count() << endl;
33 }
Tutorials -- zsh -- 80x24
kavyagupta@Kavyas-MacBook-Pro Tutorials % g++ cycle_error.cpp
kavyagupta@Kavyas-MacBook-Pro Tutorials % ./a.out
A()
B()
2
2
kavyagupta@Kavyas-MacBook-Pro Tutorials %

```

6. Why do the following three writes cause compilation errors in the C++20 compiler?

```
class Node {  
    public :  
    Node(): value(0) {}  
    const Node& foo(const Node* x) const {  
        value = 3; // Not allowed because of -----  
        x[0].value = 4; // Not allowed because of -----  
        return x[0];  
    }  
    int value;  
};  
  
int main () {  
    Node x[3], y ;  
    auto &z = y.foo(x);  
    z.value = 5; // Not allowed because of -----  
}
```

Solution: The line 'value = 3;' is not allowed since the method 'foo' is marked **const** at the end. Using **const** after the function signature and its parameters for a class method implies that the class members cannot be changed and the object itself is constant within the bounds of the function. As a result, an error is raised.

The line 'x[0].value = 4;' raises an error since the parameter 'x' is of type '**const Node***', meaning that different values can be assigned to the pointer itself, but the subscripts of the pointer cannot be assigned, since it points to constant members of class 'Node'. Likewise, '*x.value = 3;' is equally illegal. Note that '**Node const ***' also does the same thing, whereas '**Node * const**' means that the subscripts of the pointer can be assigned since it does not point to constant members of class 'Node' but the pointer itself cannot be assigned to point to a different Node or array of Nodes. '**const Node* const**' or '**Node const * const**' imply that the pointer cannot be assigned AND subscripts or dereferences cannot be assigned.

The last line 'z.value = 5;' is illegal since the datatype of 'z', as determined from the method 'foo' is '**const Node &**' NOT '**Node &**'. The implication is that 'z' refers to a Node, which is constant and hence cannot be assigned. Note that assignments to class members are as bad as assignments to the class object itself regarding constant objects in C++.

Tutorial 2

1. The span of a stock's price on i^{th} day is the maximum number of consecutive days (up to i^{th} day and including the i^{th} day) the price of the stock has been less than or equal to its price on day i .
Example: for the price sequence 2 4 6 5 9 10 of a stock, the span of prices is 1 2 3 1 5 6.
Give a linear-time algorithm that computes s_i for a given price series.

Solution: The idea here is to find out the latest (most recent) day till the i^{th} day where the price was greater than the price on day i . To do this, we will maintain a stack of indices, and for each index i , we will keep the stack in a state such that the following invariant holds for each day i :

If the stack is not empty, then the price on the day whose index is at the top of the stack is the most recent price that was strictly greater than the current price.

With this invariant, the updates to the stack and the stock span follow as:

- we remove the top index in the stack till the price corresponding to the top index becomes strictly greater, or the stack becomes empty
- if the stack is empty, the current price is the largest, so the span is the number of days so far
- otherwise the span is the number of days since the index at the top of the stack

(Such a stack is known as a “Monotonic Stack”).

The base case is that at day 1 the price is the only encountered price, and hence is the largest, so the span is 1.

The time complexity is indeed $\mathcal{O}(n)$ (n being the number of days/size of array s) because each element is pushed **exactly once** in the stack and popped **atmost once**.

Look at Algorithm 2 for the pseudo-code.

Algorithm 2: Stock Span Algorithm

Data: Array *price* of length n and an empty stack *St*

```
1 for  $i \leftarrow 0$  to  $n - 1$  do
2    $s[i] \leftarrow i + 1$  //If price[i] is  $\geq$  all the previous prices, then  $s[i]$  will be  $i + 1$ 
3   while not St.isEmpty() do
4     if  $\text{price}[\text{St.top}()] > \text{price}[i]$  then
5        $s[i] \leftarrow i - \text{St.top}()$ 
6       break
7     end
8     St.pop()
9   end
10  St.push(i)
11 end
```

2. There is a stack of dosas on a tava, of distinct radii. We want to serve the dosas of increasing radii. Only two operations are allowed:
(i) serve the top dosa
(ii) insert a spatula (flat spoon) in the middle, say after the first k , hold up this partial stack and flip it upside-down and put it back

Design a data structure to represent the tava, input a given tava, and to produce an output in sorted order. What is the time complexity of your algorithm?

This is also related to the train-shunting problem.

Solution: Read the first line of the question with the emphasis being on **stack**. We will use the array representation of a stack for the tava.

Our stack-tava has the following abstraction:

- Initialization and input: Initialize an array S of size n , to represent the stack. Using the given input, fill the array elements with the radii of dosas on the tava in the initial stack order, with the bottom dosa at index 0 and the top at $n - 1$. Maintain a variable sz which contains the current size of the stack, and initialize it to n .
- Serving the top dosa: This method is essentially carrying out the “pop” operation on our stack S . Return the top dosa in the stack, $S[sz - 1]$, and decrement sz by 1. Clearly, this method takes a constant $\mathcal{O}(1)$ time for every call.
- Flipping the top k dosas: This method is equivalent to reversing the slice of the array S from index $sz - k$ to $sz - 1$. This is simple enough: initialize an index i to $sz - k$ and another index j to $sz - 1$, swap elements at indices i and j , increment i and decrement j by 1, and repeat the swap and increment/decrement while $j > i$. This method takes $\mathcal{O}(k)$ time for every call.

Note that in this implementation, the serve operation is the only way by which the problem can be reduced in size (number of dosas reduced by 1). The flip operation, if used wisely, can give us access to dosas that can help us reduce the problem size.

With this implementation, the algorithm to serve the dosas in increasing order of the radii can be designed as follows. While the stack is not empty, repeat the following:

- Iterate over the array S and find the index of the minimum element in the array. Let this index be m . This is the position of the dosa having the smallest radius
- Flip over the top $sz - m$ dosas in the stack, using the flipping operation. This brings the smallest dosa from index m to index $sz - 1$, i.e., the top of the stack
- Serve the top dosa, using the serving operation. This pops the smallest dosa off the stack

It is easy to argue the correctness. The invariant is that we keep serving the minimum radii dosa from the remaining stack. The resulting order has to be sorted.

For computing the time complexity: in a stack of size z , finding the index of the minimum element takes $\mathcal{O}(z)$ time, flipping over the top $k \leq z$ elements takes $\mathcal{O}(k)$ time (and thus $\mathcal{O}(z)$ time) and serving off the top dosa takes $\mathcal{O}(1)$ time. Thus serving the smallest dosa takes $\mathcal{O}(z)$ time overall. This has to be repeated for all dosas, so the stack size z goes from n to 1, decrementing by 1 every time. Therefore the algorithm has a time complexity of $\mathcal{O}(n^2)$.

3. (a) Do the analysis of performance of exponential growth if the growth factor is three instead of two? Does it give us better or worse performance than doubling policy?
- (b) Can we do the similar analysis for growth factor 1.5?

Solution: Let us do the analysis for a general growth factor α . Suppose initially $N = 1$ and there are $n = \alpha^i$ consecutive pushes. So, total cost of expansion is (refer to slides for detailed explanation):

$$\begin{aligned}
 &(\alpha + 1) \cdot (\alpha^0 + \alpha^1 + \alpha^2 \dots + \alpha^{i-1}) \\
 &\quad \frac{\alpha + 1}{\alpha - 1} \cdot (\alpha^i - 1) \\
 &\quad \frac{\alpha + 1}{\alpha - 1} \cdot (n - 1)
 \end{aligned}$$

For α equals to 3 cost of expansion is $2 \cdot (n - 1)$, it's better than doubling policy. For α equals 1.5 cost of expansion is $5 \cdot (n - 1)$, which is worse than doubling policy. Trade-off involved is extra memory allocation, maximum extra memory allocated is $\alpha - 1$ times the requirement, so with increasing alpha extra memory allocation is increasing.

4. Give an algorithm to reverse a linked list. You must use only three extra pointers.

Solution: Let us initialise three pointers - *prev*, *curr*, and *next*, initialised to null, head and *head* \rightarrow *next* respectively. If head is null, then it is an empty linked list and we return. Otherwise, we will be iterating over each element in the linked list. In each iteration, perform the following updates.

- *curr* \rightarrow *next* = *prev*
- *prev* = *curr*
- *curr* = *next*
- *next* = *next* \rightarrow *next*

The loop terminates when *next* is null, in which case, set head to *curr*.

5. Give an algorithm to find the middle element of a singly linked list.

Solution: The idea is as follows:

- Make 2 pointers middle and end initialized to head of linked list
- At each step, increment middle by one and end by 2
- Continue this process until end reaches the end of linked list
- return the middle element
- return null if the head itself was null (empty linked list)

Note: When there is an odd number of elements, the above algorithm returns the element at the median location. What happens when there is an even number of elements?

6. Given two stacks *S1* and *S2* (working in the LIFO method) as black boxes, with the regular methods: "Push", "Pop", and "isEmpty", you need to implement a Queue (specifically: Enqueue and Dequeue working in the FIFO method). Assume there are n Enqueue/Dequeue operations on your queue. The time complexity of a single method Enqueue or Dequeue may be linear in n , however the total time complexity of the n operations should also be $\Theta(n)$

Solution:

Let us try to formulate a simple implementation of a queue using 2 stacks.

Label the stacks as A and B . We will use stack A to store elements of the queue sequentially, and stack B to implement the dequeue operation.

Enqueue: To insert a new element to the queue, we can just *push* the element to the top of stack A . Here, each *enqueue* operation requires $\mathcal{O}(1)$ time.

Dequeue: To remove an element from the queue in a FIFO manner, we must somehow delete the first element pushed to stack A . Since we cannot access any element of the stack other than the top one (last inserted), we will have to use stack B to implement this.

To get access to the first inserted element, we need to *pop* all the remaining elements from stack A and store them into stack B . As you can see in the figure 1, the order of elements is reversed while moving them to stack B . When we reach the first element of A , we can just *pop* it without then adding it to B .

We must now decide what to do with the elements of stack B . One idea is to simply use stack B for storing the queue in the place of stack A , but this doesn't work since the order of elements would change.

Another possible approach would be to again transfer all the elements saved in stack B to stack A and restore the queue. This would preserve the order of elements and allow new elements to be pushed to A while maintaining correctness. However, think of the worst case time complexity of this approach. For every *dequeue* operation, all the elements have to be shifted twice - once from A to B and then from B to A . This causes each *dequeue* operation to be of $\mathcal{O}(n)$ time complexity and the total time complexity of n operations to be $\mathcal{O}(n^2)$. (think when?)

The blowup in time complexity is caused due to elements being shifted back to A from B on every *dequeue*. Do we really need that? What happens if we just let the elements stay in B ?

What happens now if we have 2 successive dequeue operations? Suppose we transferred the elements of queue from stack A to B in the first operation. For the 2nd dequeue operation, we can simply *pop* from stack B . Hence we don't need to re-transfer the elements. Now, if we have n operations, our complexity reduces to $\mathcal{O}(n)$.

To sum it up: Look at Algorithm 3 for the pseudo-code

Algorithm 3: Queue using Two Stacks

```
1 Function Enqueue(element):
2   | S1.Push(element)
3 end
4 Function Dequeue():
5   | if S1.isEmpty() and S2.isEmpty() then
6   |   | throw("Queue is empty.")
7   | end
8   | if S2.isEmpty() then
9   |   | while not S1.isEmpty() do
10  |   |   | S2.Push(S1.Pop())
11  |   | end
12  | end
13  | return S2.Pop()
14 end
```

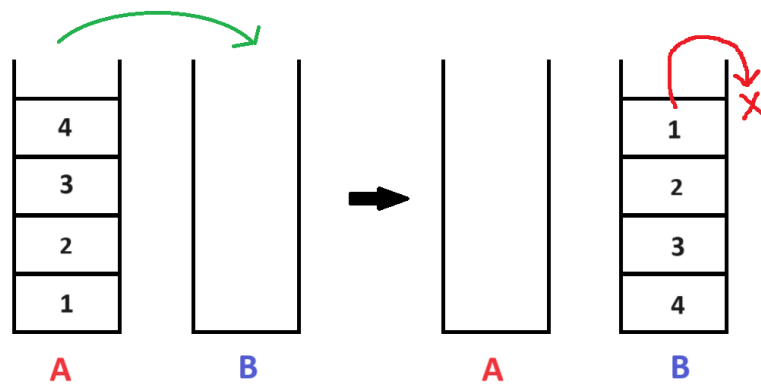


Figure 1: 2 Stacks making a Queue

Tutorial 3

1. What is the probability for the 3rd insertion to have exactly two collisions while using linear probing in the hash table?

Solution: Let the size of the hash table be $n \geq 3$. We assume that insertions take place via linear probing. Further, we assume that no deletions or insertions occur in the same sequence. Lastly, we assume that the hash function is uniform, i.e., for a random key, the likelihood of any hash value occurring is exactly $1/n$. After collisions, let the first two insertions occur at indices i and j .

If $((i+1) \bmod n) \neq j$ and $((j+1) \bmod n) \neq i$, then the third insertion cannot have two collisions. The reason is that if the hash value of the third key collides with the first key, then index $(i+1) \bmod n$ is vacant, and there is only one collision. If the hash value is at index j , then index $(j+1) \bmod n$ is vacant and hence, only one collision. Any other hash value implies no collision. Let us neglect this case.

We focus only on cases where i and j are adjacent to each other on the table. Let us re-label the indices as a and b such that $((a+1) \bmod n) = b$. In this case, for any constant $0 \leq c < n$, the absolute value of $P(a = c)$ (not conditional) can be thought of as the probability that $((i = c \text{ and } j = (i+1) \bmod n) \text{ or } (j = c \text{ and } i = (j+1) \bmod n))$. The first event occurs when the second hash value is either i or $(i+1) \bmod n$, with a chance of $(1/n) \cdot (2/n)$. The second event occurs when the second hash value is $(n+i-1) \bmod n$, with a chance of $(1/n) \cdot (1/n)$.

From the above, we can say that $P(a = c) = (3/n^2)$ for all constants c , we safely say that two collisions occur if and only if the third hash-value is a , which occurs with a probability of $\frac{1}{n}$. However, c and, therefore, a can take any value. Thus, the product $(3/n^2) \cdot (1/n)$ has to be summed up across all values (n times), giving us the required answer: $(3/n^2)$.

2. Given that k elements have to be stored using a hash function with target space n . What is the probability of the hash function having an inherent collision? What is an estimate of the probability of a collision in the insertion of n elements?

Solution: To compute the probability, we first count the negative case: no collisions at all. This means that the k keys get mapped to distinct k keys out of n . We assume that for a random key and for any constant c , $P(\text{hash}(k) = c) = 1/n$.

The first key has n choices of hash values. Given each choice of the first key, the second key has $n-1$ possible hash values such that there is no collision. Given the first two keys and hash values, the third has $n-2$ choices, and so on, till the k^{th} key has $n-k+1$ choices. Overall, the ordered k -tuple has $\prod_{i=0}^{k-1} (n-i)$ choices for no collision. All those are mutually exclusive events with likelihoods $(1/n)^k$ each, being the product of k probabilities of $1/n$ for k independent events and their respective outcomes. Thus,

$$P(\text{No collisions}) = \frac{n!}{(n-k)! \cdot n^k}$$

One minus this expression gives the likelihood of a collision.

If $k = n$, the expression boils down to $(n!/n^n)$. We apply Stirling's approximation here:

$$\sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n+1}} < n! < \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n}}$$

$$\sqrt{2\pi n} \left(\frac{1}{e}\right)^n e^{\frac{1}{12n+1}} < \frac{n!}{n^n} < \sqrt{2\pi n} \left(\frac{1}{e}\right)^n e^{\frac{1}{12n}}$$

For large n , we can neglect the terms $e^{(1/(12n+1))}$ and $e^{(1/12n)}$ as both are very close to 1. However, the term $\sqrt{2\pi n} \cdot e^{-n}$ tends to zero as n grows very large. This means that the required probability is sandwiched between terms that exponentially decay toward zero.

The implication of the above is the likelihood that all n keys get uniquely mapped to the n possible hash values decay exponentially with n , and thus, collisions are almost certain for large n . This can be linked with the birthday paradox in that in a class of 365; it is almost impossible for all students to have different birthdays; even a single tutorial batch has a high likelihood of at least two birthdays being the same. In the context of hashing, although the expected number of terms sharing a hash value is one per hash value, we expect a good number of collisions to occur during hashing.

3. Let $C(i)$ be the chain of array indices that are queried to look for a key k in linear probing where $h(k) = i$.
- (a) How does this chain extend by an insertion, and how does it change by a deletion?
 - (b) A search for a key k ends when an empty cell is encountered. What if we mark the end of $C(i)$ with an end marker? We stop the search when this marker is encountered. Would this work? Would this be efficient?
 - (c) Is there a way of not using tombstones?

Solution: To recall, linear probing covers indices $i, i + 1, \dots$ until an empty index is found (then k not found in the table) or the cell stores k . Note that tombstones do not count as empty cells in the hash table; counting them so after deletion would break existing chains and result in false negatives when searching. Including them means no change in these chains, answering part of the first (a) subproblem.

If the chain terminates at k (key found), then inserting a new element does not change anything. If the added element is k itself, then it occupies the position at the end of the chain, a vacant cell. So, no change. In the other case, the chain extends by one if the inserted element is at the end of the chain and is not k . The chain could extend a lot more if the inserted element joins the existing chain with another series of filled cells in the table. Note that the chain, except for the end, has only filled or tombstone cells; thus, the insertion cannot occur in the middle of the chain. If the insertion of the element that is not k does not occur at the end of the previous chain for searching k , then there is no change.

The answer to the second (b) part is that it would work but not be any more efficient than the standard algorithm. For this to work, an empty cell once encountered during the search should be marked an 'end.' Otherwise, the first search would go on an infinite loop, looking for an end marker. If this is the case, the search must look for both end marks and empty cells in each iteration, making it slightly less efficient but at the same time complexity.

For the third (c) part, a slight modification should be made to the algorithm discussed above: all cells should be marked 'end' at the start. Otherwise, a deletion without a tombstone would cause false negatives in future searches, just like the discussions without tombstones. If this modification is made, then we basically replace "empty" with an "end" marker and "tombstone" marker with "empty". The algorithm is, otherwise, identical to the normal algorithm with tombstones, and hence, we really do not avoid tombstones, although we do not use them. Note that marking the freed cells as 'end' on deletions would be as bad as not using tombstones, which results in false negatives.

4. Let $m = 11$, $h_1(k) = (k \bmod 11)$, $h_2(k) = 6 - (k \bmod 6)$.
Let us use the following hash function for an open addressing scheme.

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$$

1. What will be the state of the table after insertions of 41, 22, 44, 59, 32, 31, and 74?
2. Let $h_2(k) = p - (k \bmod p)$. What should be the relationship between p and m such that h is a valid function for linear probing?
3. What is the average number of probes for an unsuccessful search if the table has α load factor?
4. What is the average time for a successful search?

Solution: To recall, the standard linear probing uses the function $h(k, i) = (h(k) + i) \bmod m$ where m is the size of the hash table [Here 11]. The way it works is that we check the cells indexed at $h(k, 0), h(k, 1), h(k, 2), \dots$ till we find an empty cell $h(k, a)$ for some a . The same logic should be extended here, except $h(k, i)$, which is defined in the problem statement.

Let the table before adding 41 be

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

Now, $h_1(41) = 8$ and $h_2(41) = 1$. This means that $h(41, 0) = 8$ and cell 8 is empty. So, we add 41 there:

								41		
--	--	--	--	--	--	--	--	----	--	--

$h_1(22) = 0$ and cell 0 is empty so we can add 22 there. [$h(k, 0) = h_1(k)$]

22								41		
----	--	--	--	--	--	--	--	----	--	--

$h_1(44) = 0$ but cell 0 is filled. Now, $h_2(44) = 4$ and $h(44, 1) = 4$ which is an empty tile.

22				44				41		
----	--	--	--	----	--	--	--	----	--	--

$h_1(59) = 4$ which is a filled cell. $h_2(59) = 1$ and $h(59, 1) = 5$ is an empty tile.

22				44	59			41		
----	--	--	--	----	----	--	--	----	--	--

$h_1(32) = 10$ and $h_1(31) = 9$ are both empty and different cells, so we can add them.

22				44	59			41	31	32
----	--	--	--	----	----	--	--	----	----	----

Lastly, $h_1(74) = 8$, which is occupied. $h_2(74) = 4$. Now, $h(74, 1) = 1$ which is free.

22	74			44	59			41	31	32
----	----	--	--	----	----	--	--	----	----	----

This is how the table looks after inserting the elements in the given order.

Moving on to the next part, a valid function for linear probing should eventually cover all cells in the grid as i increments. This is not possible if $\gcd(h_2(k), m) > 1$, for one or more values of k . The reason is that $h(k, i) \bmod \gcd(h_2(k), m)$ will be fixed, and hence, $h(k, i) \bmod m$ will only take some values. But if that gcd is 1, then it means that for all $0 \leq i < m$, $\gcd \cdot i \bmod m$ will have a unique value, and the number of these unique values will be m , so all cells of the table will be covered.

Note that $\gcd(h_2(k), m) = 1$ should hold for all values of k , not just for any one k . Keeping $h_2(k) = p - (k \bmod p)$ can take any value in $1, 2, \dots, p$ depending on the value of k . Thus, none of the first p natural numbers should share a common factor with m . In other words, the smallest prime factor (not 1) of m , the size of the hash table, should be larger than p . It is not necessary that m is prime – for $p = 6$, $m = 7 \cdot 13 = 91$ also works.

Exact expression of expected number of probes in a search (unsuccessful/successful) is mathematically complex hence we'll prove an appropriate upper bound for them. Claim: The mentioned expression for an unsuccessful one is atmost $1/(1 - \alpha)$. This is assuming $\alpha < 1$.

Proof: Let X be the number of probes made in an unsuccessful search, and let A_i be the event that an i^{th} probe occurs and it is to an occupied slot. Then $X \geq i$ when the event $A_1 \cap A_2 \cap \dots \cap A_{i-1}$ occurs. Hence,

$$\begin{aligned} P(X \geq i) &= P(A_1 \cap A_2 \cap \dots \cap A_{i-1}) \\ &= P(A_1) \times P(A_2|A_1) \times P(A_3|A_2 \cap A_1) \times \dots \times P(A_{i-1}|A_1 \cap \dots \cap A_{i-2}) \end{aligned}$$

$P(A_1) = n/m$ because there are n elements and m slots. For $j > 1$, the probability that there is a j^{th} probe and it is to an occupied slot, given that the first $j - 1$ probes were to occupied slots, is $(n - j + 1)/(m - j + 1)$ because we are finding one of the remaining $(n - (j - 1))$ elements in one of the $(m - (j - 1))$ unexamined slots, and we have uniform hashing so each of the slots is equally likely to be chosen next.

Because we have $n < m$, we have $(n - j)/(m - j) \leq n/m \forall 0 \leq j < m$. Thus,

$$P(X \geq i) = \frac{n}{m} \times \frac{n-1}{m-1} \times \dots \times \frac{n-i+2}{m-i+2} \leq \left(\frac{n}{m}\right)^{i-1} = \alpha^{i-1}$$

So $E[X] = \sum_{i=1}^{m-1} P(i^{th} \text{ probe occurs}) = \sum_{i=1}^{m-1} P(X \geq i) \leq \sum_{i=1}^{m-1} \alpha^{i-1} \leq \sum_{i=1}^{\infty} \alpha^{i-1} = 1/(1 - \alpha)$. Hence proved.

Claim: Expected number of probes for a successful search is atmost $\frac{1}{\alpha} \ln(\frac{1}{1-\alpha})$.

Proof: A search for a key k reproduces the same probe sequence as when the element with key k was inserted. If k was the $(i + 1)^{th}$ key inserted into the hash table, there were only i keys in the table at the time, so the expected number of probes made in a search for k is at most $1/(1 - i/m) = m/(m - i)$ (as an insertion is equivalent to an unsuccessful search, here $\alpha = i/m$). Averaging over all n keys in the hash table, we get,

$$\begin{aligned} \frac{1}{n} \sum_{i=0}^{n-1} \frac{m}{m-i} &= \frac{m}{n} \sum_{i=0}^{n-1} \frac{1}{m-i} \\ &= \frac{1}{\alpha} \sum_{k=m-n+1}^m \frac{1}{k} \\ &\leq \frac{1}{\alpha} \int_{m-n}^n \frac{1}{x} dx \\ &= \frac{1}{\alpha} \ln\left(\frac{m}{m-n}\right) \\ &= \frac{1}{\alpha} \ln\left(\frac{1}{1-\alpha}\right) \end{aligned}$$

Hence Proved.

5. Suppose you want to store a large set of key-value pairs, for example, (name, address). You have operations in this set: addition, deletion, and search of elements. You also have queries about whether a particular name or address is in the set, and if so, then count them and delete all such entries. How would you design your hash tables?

Solution: Let us recall that in the classic hashing, the hash function operates on the key, and the value is simply stored alongside the key. An intuitive idea is to keep the key (from the point of view of hashing) as **both** the name and the address, with no special value otherwise. The hash function operates on this and gives us the index. Unfortunately, this gives rise to a problem: there are multiple different hash values for a given name, and hence, searching for the existence of an entry given a name (but not address – we do not know which addresses are there given this name) is as bad as linearly iterating through the table.

One way to get around this is to use a two-level hash table: the outer table N stores the keys as the names alone, and each value corresponding to a name points to a hash table itself. This inner hash table $N[name]$ stores only addresses and only addresses where the name is the corresponding key of the outer table. This makes finding elements based on names easy but does not solve the problem of searching for given addresses.

To solve that, we can use a double two-level hash table: one two-level table N as described above and one two-level table A done the other way round, i.e., A stores the keys as addresses and values as pointers to inner hash tables $A[address]$ that store names corresponding to addresses. Note that this requires additional storage space (double) and more maintenance.

Specifically, it is important to ensure integrity: for every $(name, address)$ pair in the set, $N[name]$ should exist and contain $address$, and $A[address]$ should exist and contain $name$. Insertion should be done in both tables at once. On the other hand, if any $(name, address)$ pair does not exist, then $A[address]$, if it exists, should not contain $name$, and $N[name]$, if it exists, should not contain $address$. This should be reflected in the deletion of pairs.