

## Chapter 9

# Perfect Hashing and Graphs.

In the past week, we have seen how to store keys in hash tables so that the number of mapped keys in a specific cell is  $O(1)$  in expectation. The table is constructed using a hashing function  $h : \text{key space} \rightarrow m\text{cells}$ , randomly chosen from a universal hash function family. This function maps keys to cells, and in each cell, the keys are stored using a linked list. The cost of supported subroutines depends on the length of the list. We named any pair of different keys  $x \neq y$  that are mapped to the same cell in the table, namely  $h(x) = h(y)$ , a collision.

Perfect hashing is a method to ensure that no collision occurs, it works only if all keys are given in advance and they are unique, meaning that the table doesn't support insertion. The idea is as follows, we sample an hash function, and then check if, for all  $x, y$  in the input, it holds that  $h(x) \neq h(y)$ . If so then we continue. Otherwise we repeat.

```
1 let collision  $\leftarrow$  True
2 while collision do
3   collision  $\leftarrow$  False
4   let  $T$  be array at length  $m$ 
5    $h \leftarrow$  sample uniformly random from universal hash family  $\mathcal{H}$ 
6   for  $x \in x_1, x_2, \dots, x_n$  do
7     if  $T_{h(x)}$  is not empty then
8       collision  $\leftarrow$  True
9       break the for-loop
10    end
11    else
12       $T_{h(x)} \leftarrow x$ 
13    end
14  end
15 end
16 return  $T, h$ 
```

**Algorithm 1:** perfect-hashing( $x_1, x_2, \dots, x_n$ )

**Question.** What is the probability of choosing  $h$  with no collisions on the first trial? Notice that the answer depends on  $m$ . (To see this, imagine the case where  $m = 1$ . In this case, there must be collisions.) Therefore, the correct question is:

for what values of  $m$  do we succeed in finding a hash function with no collisions on the first trial? Let  $X_{x,y}$  be the indicator of the event  $h(x) = h(y)$ . The expected number of collisions is then:

$$\mathbf{E} \left[ \sum_{x \neq y} X_{x,y} \right] = \sum_{x \neq y} \mathbf{E} [X_{x,y}] = \binom{n}{2} \frac{1}{m}$$

Now, we would like to answer for what value of  $m$  there is no collision. Therefore, if we take  $m = n^2$ , then the expected number of collisions is less than  $1/2$ . By the Markov inequality, the probability of having more than one collision is less than:

$$P \left( \sum_{x \neq y} X_{x,y} > 1 \right) \leq \mathbf{E} \left[ \sum_{x \neq y} X_{x,y} \right] = \frac{1}{2}$$

And therefore the expected number of rounds is less than:

$$\mathbf{E} [\text{rounds}] = \sum_{t=0}^{\infty} t P(t \text{ rounds}) \leq \sum_{t=0}^{\infty} t \frac{1}{2^{t-1}} = O(1)$$

**Question.** What is the space complexities? We have to allocate an array of length  $m$  which is  $\Theta(n^2)$  memory. Is that good? So remember that in standard hash tables, the expected number of elements that were hashed into the same cell as the key  $x$  is

$$1 + \frac{n-1}{m}$$

Taking  $m = \Theta(n)$  is enough to ensure that the expected running time of insertion/deletion/access is  $O(1)$ . This raises the question of whether the space complexity of perfect hashing can be reduced to linear.

### 9.0.1 Perfect Hashing in Linear Space.

The idea is as follows: we will use a two-stage hashing process. In the first stage, keys will be mapped to hash tables instead of cells. Each hash table will be constructed using perfect hashing and may require a space that is quadratic in the number of elements stored in it (which were mapped to it in the first stage). Therefore, if we denote by  $n_i$  the number of elements mapped to the  $i$ th hash table, the space cost will be  $\sum_i n_i^2$ . Instead of starting over when a collision occurs, we will do so when  $\sum_i n_i^2 > 4n$ . So, now it's left to show that we expect  $\sum_i n_i^2$  to be linear, which implies that the expected number of rounds is constant.

$$n_i^2 = 2 \binom{n_i}{2} + n_i$$

On the other hand,  $\sum_i \binom{n_i}{2}$  is exactly the number of collisions, as for any  $i$ ,  $\binom{n_i}{2}$  counts the number of distinct pairs in the  $i$ th table, which is equivalent to counting

```

1 let toomanycollisions  $\leftarrow$  True
2 while toomanycollisions do
3   toomanycollisions  $\leftarrow$  False
4   let  $T$  be array at length  $m$ 
5   initialize any  $T_i$  to be an empty linked list.
6    $h \leftarrow$  sample uniformly random from universal hash family  $\mathcal{H}$ 
7   for  $x \in x_1, x_2, \dots, x_n$  do
8      $T_{h(x)}.insert(x)$ 
9      $T_{h(x)}.size = T_{h(x)}.size + 1$ 
10  end
11  if  $\sum_i T_{h(i)}.size^2 \geq \mu$  then
12    toomanycollisions  $\leftarrow$  True
13  end
14 end
15 let  $H$  be an array at length  $m$ 
16 for  $i \in [m]$  do
17    $T_i, h_i \leftarrow$  hash the elements in  $T_i$  using
18   perfect hashing.
19 end
20 return  $T, h$ 

```

**Algorithm 2:** perfect-hashing-linear-space( $x_1, x_2, \dots, x_n$ )

the number of  $x \neq y$  such that  $h(x) = h(y) = i$ . Thus,

$$\begin{aligned}
 \mathbf{E} \left[ \sum_i n_i^2 \right] &= \mathbf{E} \left[ \sum_i 2 \binom{n_i}{2} + n_i \right] = 2 \cdot \mathbf{E}[\text{collisions}] + \mathbf{E} \left[ \sum_i \overbrace{n_i}^n \right] \\
 &= 2 \cdot \binom{n}{2} \frac{1}{m} + n
 \end{aligned}$$

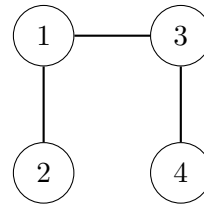
Therefore, by choosing  $m = 4n$  for the first stage, the probability of failing to choose a proper hash function is less than  $1/2$ .

## 9.1 Graphs

### 9.1.1 Definitions, Examples, and Basics

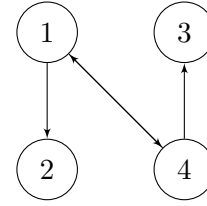
**Definition 9.1.1.** A **non-directed graph**  $G$  is a pair of two sets -  $G = (V, E)$  -  $V$  being a set of vertices and  $E$  being a set of couples of vertices, which represent edges ("links") between vertices.

**Example:**  $G = (\{1, 2, 3, 4\}, \{\{1, 2\}, \{1, 3\}, \{3, 4\}\})$  is the following graph:



**Definition 9.1.2.** A **directed graph**  $G$  is a pair of two sets -  $G = (V, E)$  -  $V$  being a set of vertices and  $E \subseteq V \times V$  being a set of directed edges ("arrows") between vertices.

**Example:**  $G = (\{1, 2, 3, 4\}, \{(1, 2), (1, 4), (4, 1), (4, 3)\})$  is the following graph (note that it has arrows):



**Definition 9.1.3.** A **weighted graph** composed by a graph  $G = (V, E)$  (either non-directed or directed) and a weight function  $w : E \rightarrow \mathbb{R}$ . Usually (but not necessary), we will think about the quantity  $w(e)$ , where  $e \in E$ , as the length of the edge.

Now that we see graphs with our eyes, we can imagine all sorts of uses for them... For example, they can represent the structure of the connections between friends on Facebook, or they can even represent which rooms in your house have doors between them.

**Remark 9.1.1.** Note that directed graphs are a **generalization** of non-directed graphs, in the sense that every non-directed graph can be represented as a directed graph. Simply take every non-directed edge  $\{v, u\}$  and turn it into two directed edges  $(v, u), (u, v)$ .

**Remark 9.1.2.** Note that most of the data structures we discussed so far - Stack, Queue, Heap, BST - can all be implemented using graphs.

Now let's define some things in graphs:

**Definition 9.1.4.** (Path, circle, degree)

1. A **simple path** in the graph  $G$  is a series of unique vertices (that is, no vertex appears twice in the series)  $v_1, v_2, \dots, v_n$  that are connected with edges in that order.
2. A **simple circle** in the graph  $G$  is a simple path such that  $v_1 = v_n$ .
3. The **distance** between two vertices  $v, u \in V$  is the length of the shortest path between them ( $\infty$  if there is no such path).

**Remark 9.1.3.** Note that for all  $u, v, w \in V$  the triangle inequality holds regarding path lengths. That is:

$$\text{dist}(u, w) \leq \text{dist}(u, v) + \text{dist}(v, w)$$

**Definition 9.1.5.** (connectivity)

1. Let  $G = (V, E)$  be a non-directed graph. A **connected component** of  $G$  is a subset  $U \subseteq V$  of maximal size in which there exists a path between every two vertices.
2. A non-directed graph  $G$  is said to be a **connected graph** if it only has one connected component.

3. Let  $G = (V, E)$  be a directed graph. A **strongly connected component** of  $G$  is a subset  $U \subseteq V$  of maximal size in which for any pair of vertices  $u, v \in U$  there exist both directed path from  $u$  to  $v$  and a directed path from  $v$  to  $u$ .

**Claim 9.1.1.** Let  $G = (V, E)$  be some graph. If  $G$  is connected, then  $|E| \geq |V| - 1$

*Proof.* We will perform the following process: Let  $\{e_1, \dots, e_m\}$  be an enumeration of  $E$ , and let  $G_0 = (V, \emptyset)$ . We will build the graphs  $G_1, G_2, \dots, G_m = G$  by adding edges one by one. Formally, we define -

$$\forall i \in [m] \quad G_i = (V, \{e_1, \dots, e_i\})$$

$G_0$  has exactly  $|V|$  connected components, as it has no edges at all. Then  $G_1$  has  $|V| - 1$ . From there on, any edges do one of the following:

1. Keeps the number of connected components the same (the edge closes a cycle)'
2. Lowers the number of connected components by 1 (the edges does not close a cycle)

So in general, the number of connected components of  $G_i$  is  $\geq |V| - i$ . Now, if  $G_m = G$  is connected, it has just one connected component! This means:

$$1 \geq |V| - |E| \implies |E| \geq |V| - 1$$

□

### 9.1.2 Graph Representation

Okay, so now we know what graphs are. But how can we represent them in a computer? There are two main options. The first one is by **array of adjacency lists**. Given some graph  $G$ , every slot in the array will contain a linked list. Each linked list will represent a list of some node's neighbors. The second option is to store edges in an **adjacency matrix**, a  $|V| \times |V|$  binary matrix in which the  $v, u$ -cell equals 1 if there is an edge connecting  $u$  to  $v$ . That matrix is denoted by  $A_G$  in the example below. Note that the running time analysis might depend on the underline representation.

**Question.** What is the memory cost of each of the representations? Note that while holding an adjacency matrix requires storing  $|V|^2$  bits regardless of the size of  $E$ , Maintaining the edges by adjacency lists costs linear memory at the number of edges and, therefore, only  $\Theta(|V| + |E|)$  bits.

**Example.** Consider the following directed graph:

### 9.1.3 Breadth First Search (BFS)

One natural thing we might want to do is to travel around inside a graph. That is, we would like to visit all of the vertices in a graph in some order that relates to the edges. Breadth-first search constructs a breadth-first tree, initially containing only its root, which is the source vertex  $s$ . Whenever the search discovers a white

Figure 9.3: Presenting  $G$  by array of adjacency lists.Figure 9.4: Presenting  $G$  by adjacency matrix.

vertex  $v$  in the course of scanning the adjacency list of a gray vertex  $u$ , the vertex  $v$  and the edge  $(u, v)$  are added to the tree. We say that  $u$  is the predecessor or parent of  $v$  in the breadth-first tree. Since every vertex reachable from  $s$  is discovered at most once, each vertex reachable from  $s$  has exactly one parent. (There is one exception: because  $s$  is the root of the breadth-first tree, it has no parent.) Ancestor and descendant relationships in the breadth-first tree are defined relative to the root  $s$  as usual: if  $u$  is on the simple path in the tree from the root  $s$  to vertex  $v$ , then  $u$  is an ancestor of  $v$ , and  $v$  is a descendant of  $u$ . The breadth-first-search procedure BFS on the following page assumes that the graph  $G = (V, E)$  is represented using adjacency lists. It denotes the queue by  $Q$ , and it attaches three additional attributes to each vertex  $v$  in the graph:

1.  $v.\text{visited}$  is a boolean flag which indicate wheter  $v$  was allready visited.
2.  $\pi(v)$  is  $v$ 's predecessor in the breadth-first tree. If  $v$  has no predecessor because it is the source vertex or is undiscovered, then  $\pi(v)$  is None/NULL.

Correctness: The example should be enough to explain the correctness. A concrete proof can be found in the book, page 597.

Runtime: We can analyse the runtime line-by-line:

- Lines 1-2:  $|V|$  operations, all in  $O(1)$  runtime, for a total of  $O(|V|)$ .
- Lines 3-6:  $O(1)$
- Lines 7-8: First we need to understand the number of times the *while* loop iterates. We can see that every vertex can only enter the queue ONCE (since it is then tagged as "visited"), and therefore it runs  $\leq |V|$  times. All operations are  $O(1)$ , and we get a total of  $O(|V|)$ .

```

1 for  $v \in V$  do
2    $v.visited \leftarrow \text{False}$ 
3 end
4  $Q \leftarrow \text{new Queue}$ 
5  $Q.Enqueue(s)$ 
6  $s.visited \leftarrow \text{True}$ 
7 while  $Q$  is not empty do
8    $u \leftarrow Q.Dequeue()$ 
9   for neighbor  $w$  of  $u$  do
10    if  $w.visited$  is False then
11       $w.visited \leftarrow \text{True}$ 
12       $\pi(w) \leftarrow u$ 
13       $Q.Enqueue(w)$ 
14    end
15  end
16 end

```

**Algorithm 3:** BFS( $G, s$ )

- Lines 9-13: Next, we want to understand the number of times this *for* loop iterates. The for loop starts iterating once per vertex, and then the number of its iterations is the same as the number of neighbors that this vertex has. Thus, it runs  $O(|E|)$  times.

So all in all we get a runtime of  $O(|V| + |E|)$

#### 9.1.4 Usage of BFS

Now we have a way to travel through a graph using the edges. How else can we use it?

**Exercise:** Present and analyse an algorithm  $CC(G)$  which receives some undirected graph  $G$  and outputs the number of connected components in  $G$  in  $O(|V| + |E|)$ .

**Solution:** Consider the following algorithm:

```

1 count  $\leftarrow 0$ 
2 for  $v \in V$  do
3   if  $v.visited = \text{False}$  then
4     count  $\leftarrow$  count + 1
5     BFS( $G, v$ )
6   end
7 end
8 return count

```

**Algorithm 4:** CC( $G$ )