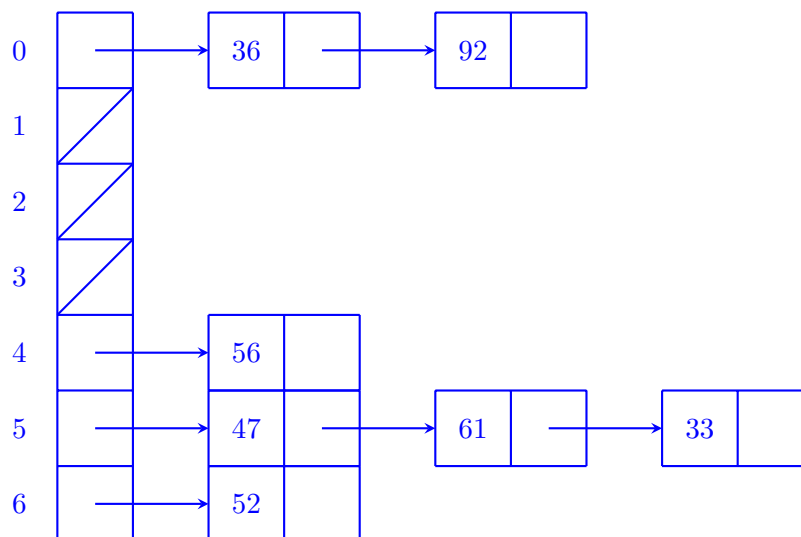# Problem Set 3

**Name:** Your Name

**Collaborators:** Name1, Name2

**Problem 3-1.**

(a) Known: $h(k) = (10k + 4) \bmod 7$
$h(47) = 5, h(61) = 5, h(36) = 0, h(52) = 6, h(56) = 4, h(33) = 5, h(92) = 0$

**(b)** Known: $h(k) = ((10k + 4) \bmod c) \bmod 7$

```python
def main():
    data = [47, 61, 36, 52, 56, 33, 92]
    'In python, set and dictionary are both based on hash table.'
    record = set()
    for c in range(1, 100):
        flag = True
        for a in data:
            x = ((10 * a + 4) % c) % 7
            if x in record:
                flag = False
                break
            record.add(x)
        if flag:
            print(c)
            break
        record.clear()

if __name__ == '__main__':
    main()
```

Answer is 13, then there is no collision.

**Problem 3-2.**

**(a)** Known: $H = \{h_{ab}(k) = (ak + b) \bmod n \mid a, b \in \{0, ..., n-1\} \text{ and } a \neq 0\}$

$\exists k_1, k_2, h(k1) = h(k2) \Leftrightarrow ak_1 + b \equiv ak_2 + b \pmod{n} \Leftrightarrow a(k_1 - k_2) \equiv 0 \pmod{n}, \forall a \in \mathbb{Z} \Leftrightarrow n \mid (k_1 - k_2)$

**(b)** $n \mid (\lfloor \frac{k_1 n}{u} \rfloor) - \lfloor \frac{k_2 n}{u} \rfloor$

But $k_1, k_2 < u \Rightarrow \lfloor \frac{k_{1(2)} n}{u} \rfloor < n$, we can deduce that the difference must be 0: $\lfloor \frac{k_1 n}{u} \rfloor = \lfloor \frac{k_2 n}{u} \rfloor$.

We can take $k_1 = 1, k_2 = 2$ because then $\frac{u}{k_{1(2)}} \gg n \Rightarrow \lfloor \frac{k_1 n}{u} \rfloor = \lfloor \frac{k_2 n}{u} \rfloor = 0$

**(c)** Seen the cour: the probability maximum is $\frac{1}{m}$

## Problem 3-3.

(a) Use **radix sort**: a string is divided into multiple ASCII characters, which can be described by a number from 0 to 127. The length of string is fixed. So we can get the time: $\Theta(nlog_4n)$. In a comparative model, each comparison of string will cost extra $\Theta(log_4n)$ time, which leads to the time total: $O(nlog^2n)$

(b)  1. $n \gg 800,000$, then $800,000$ is not a big number, so we can use **count sort**. The space cost will depend on $n$. The time cost $\Theta(n + 800,000) \simeq \Theta(n)$

  2. Otherwise, we can use merge sort et etc.

(c) Multiply by $n^3$, then use radix sort: $\Theta(n)$

(d) This is a comparative model, where we can use merge sort: $\Theta(nlogn)$

## Problem 3-4.

(a) Use a hash table to store $(i, r - b_i)$ pair.

(b)

## Problem 3-5.

(a) Use hash table. For a string A, we extract a substring of which length is k. We sort it by using bucket sort: $\Theta(k + 26) \simeq \Theta(k)$. The result of sort is the key of hash table. The value is a index sort.

When we search the anagram substring count of B, we just ust bucket sort to sort B at first and then return the size of $h(key)$

(b) Use the method of problem above is fine.

```
1   def count_anagram_substrings(T, S):
2       '''
3       Input:   T | String
4                S | Tuple of strings S_i of equal length k < |T|
5       Output: A | Tuple of integers a_i:
6                  | the anagram substring count of S_i in T
7       '''
8       A = []
9       ##################
10      # YOUR CODE HERE #
11      ##################
12      k = len(S[0])
13      record = {}
14      for i in range(0, len(T) - k + 1):
15          R = bucket_sort(T[i:i + k])
16          if R not in record:
17              record[R] = 1
18          else:
```

```
19            record[R] += 1
20      for s in S:
21          R = bucket_sort(s)
22          if R in record:
23              A.append(record[R])
24          else:
25              A.append(0)
26      return tuple(A)
27
28  def bucket_sort(T):
29      A = [0] * 26
30      for t in T:
31          A[ord(t) - 97] += 1
32      S = ""
33      for i in range(26):
34          S += chr(i + 97) * A[i]
35      return S
```

(c) Submit your implementation to `alg.mit.edu`.