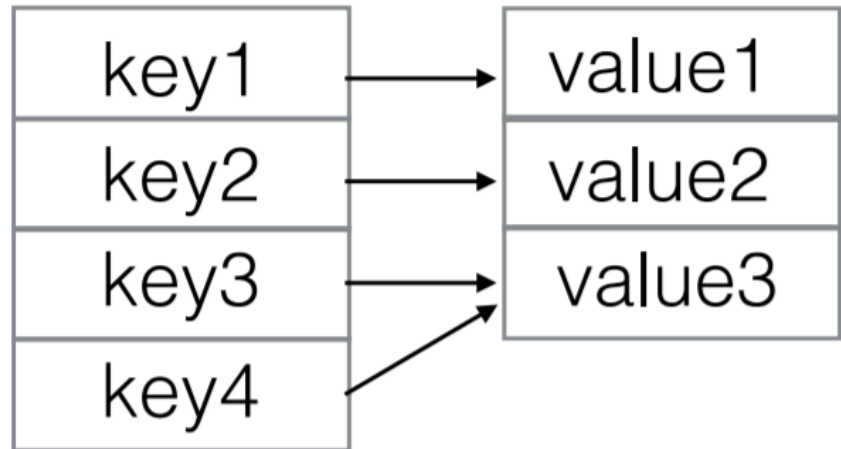


# Hashing

# Map ADT

- A *map* is a collection of *(key, value)* pairs
- Keys are unique, values may not be
- Two operations:
  - `get(key)`
  - `put(key, value)`

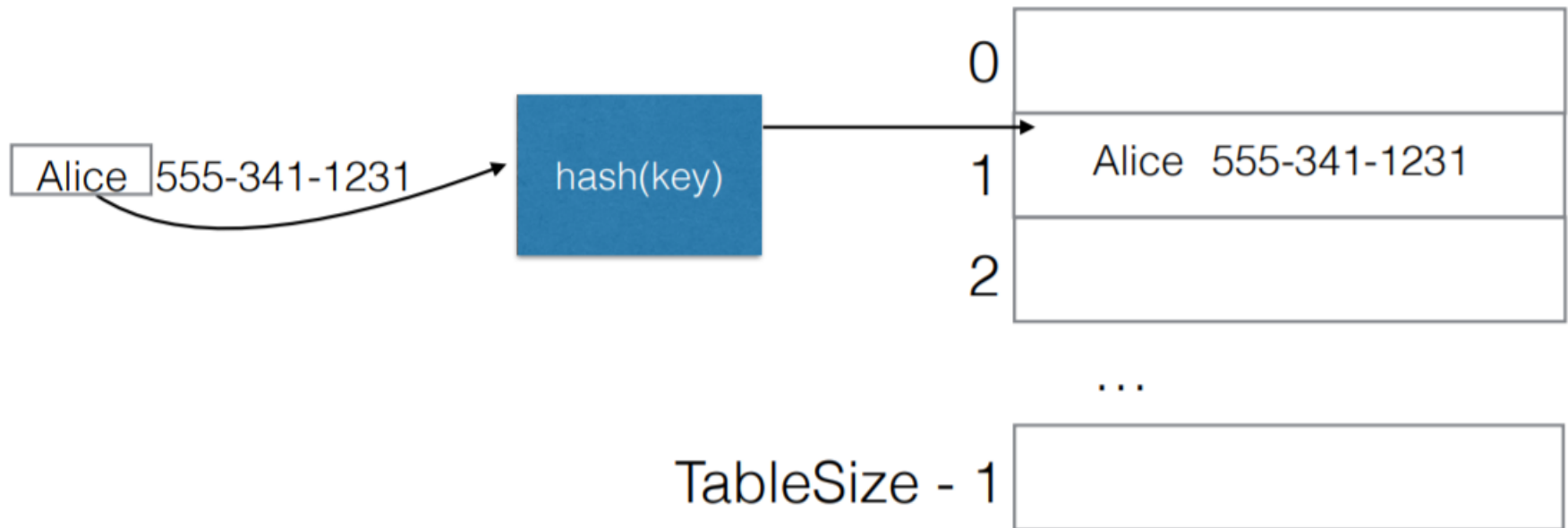


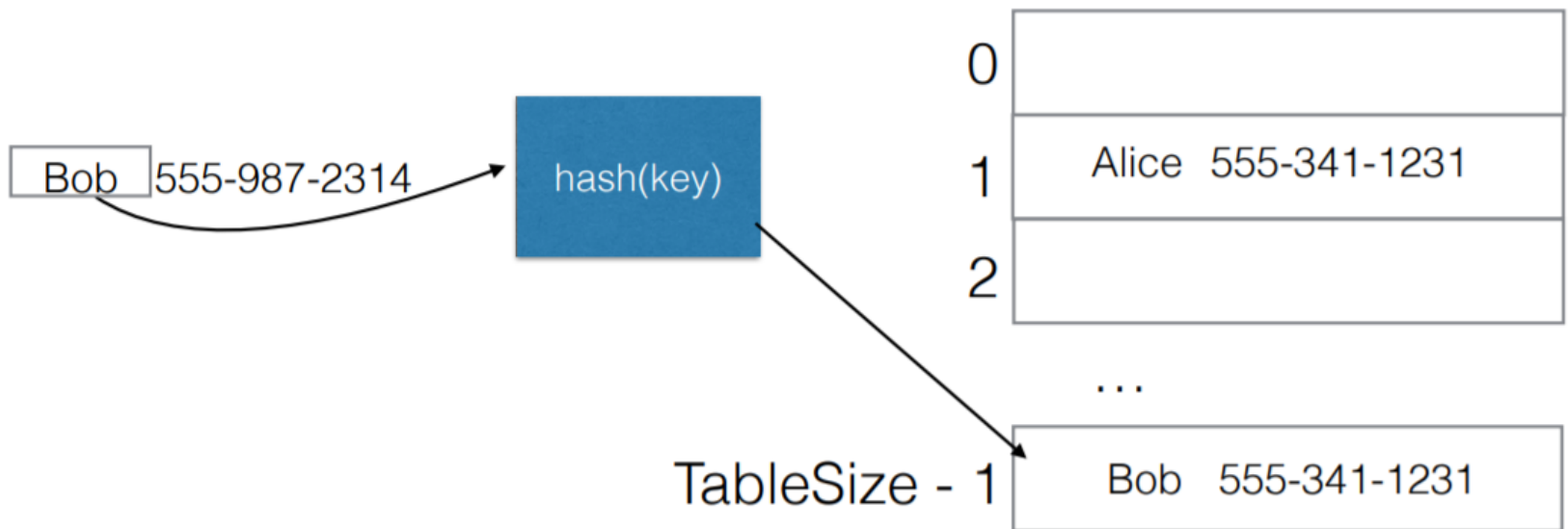
# Implementing Maps

- Store *(key, value)* pairs in a Set
  - For example, store them in an AVL Tree, with the comparator comparing keys
  - $O(\log(n))$  get and put
- Use an array of values
  - Only integer keys permitted, the array may need to be large
  - $O(1)$  get and put
- Use hash tables
  - An extension of the array idea, but without needing very large arrays

# Hash Tables

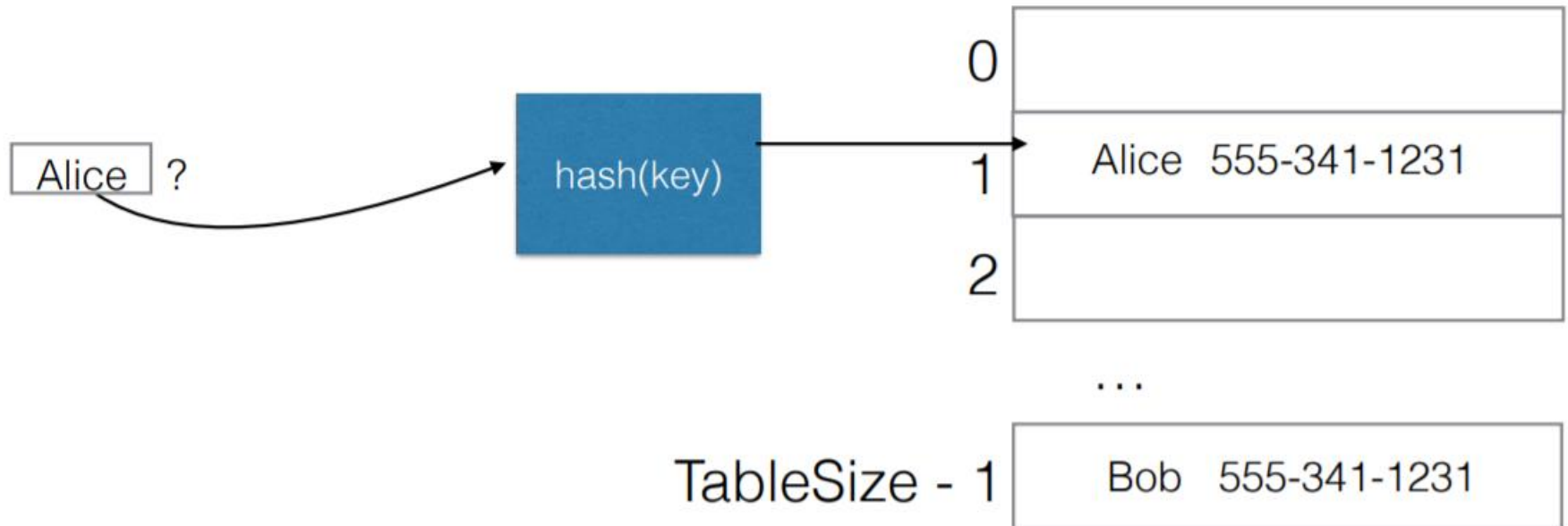
- Define a table (array) of length *TableSize*
- Define a function `hash(key)` that maps keys to integer indices in the range  $0 \dots \text{TableSize} - 1$





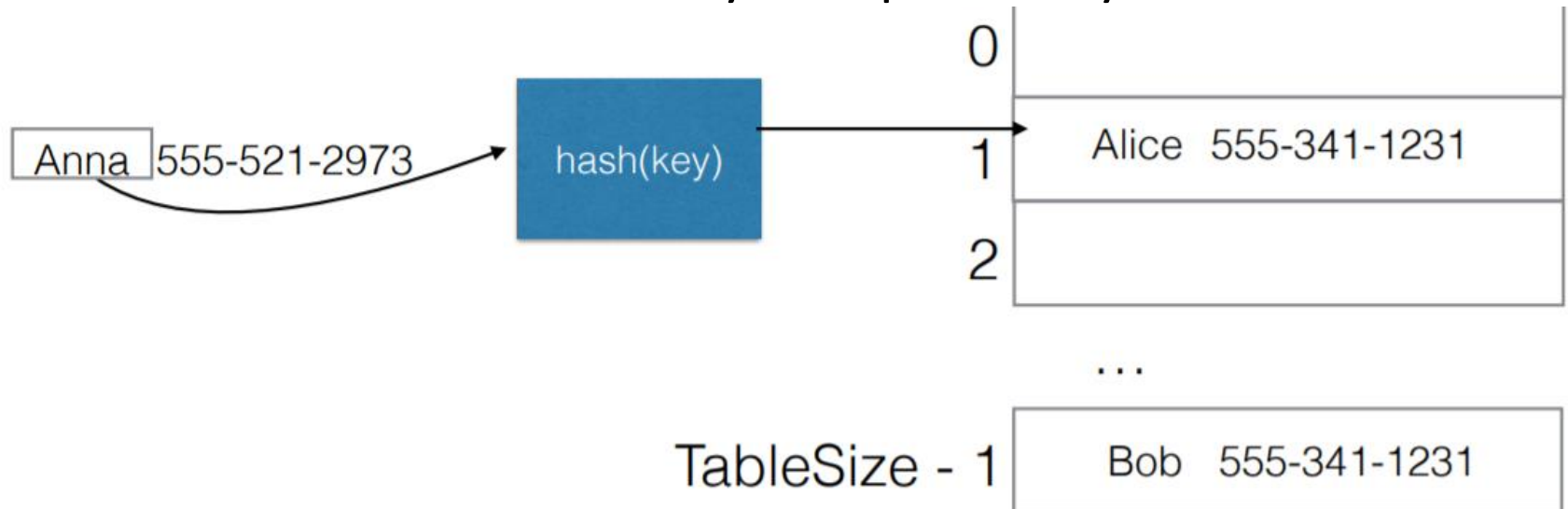
# Hash Tables

- Assuming  $\text{hash}(\text{key})$  takes constant time, get and put run in  $O(1)$



# Hash Tables Collisions

- Problem: there can be an infinite number of keys, but only *TableSize* entries in the array
  - Want to use a hash function that distributes items in the array evenly
  - Need to deal with *collisions* – situations where a new item hashes to an already-occupied array cell



# Choosing a Hash Function

- Need to
  - Spread out the keys as much as possible in the table
    - If the function maps a lot of the keys to the same area of the table, there will be collisions
  - Make use of all table cells
    - Otherwise we are wasting space



# Choosing a Hash Function

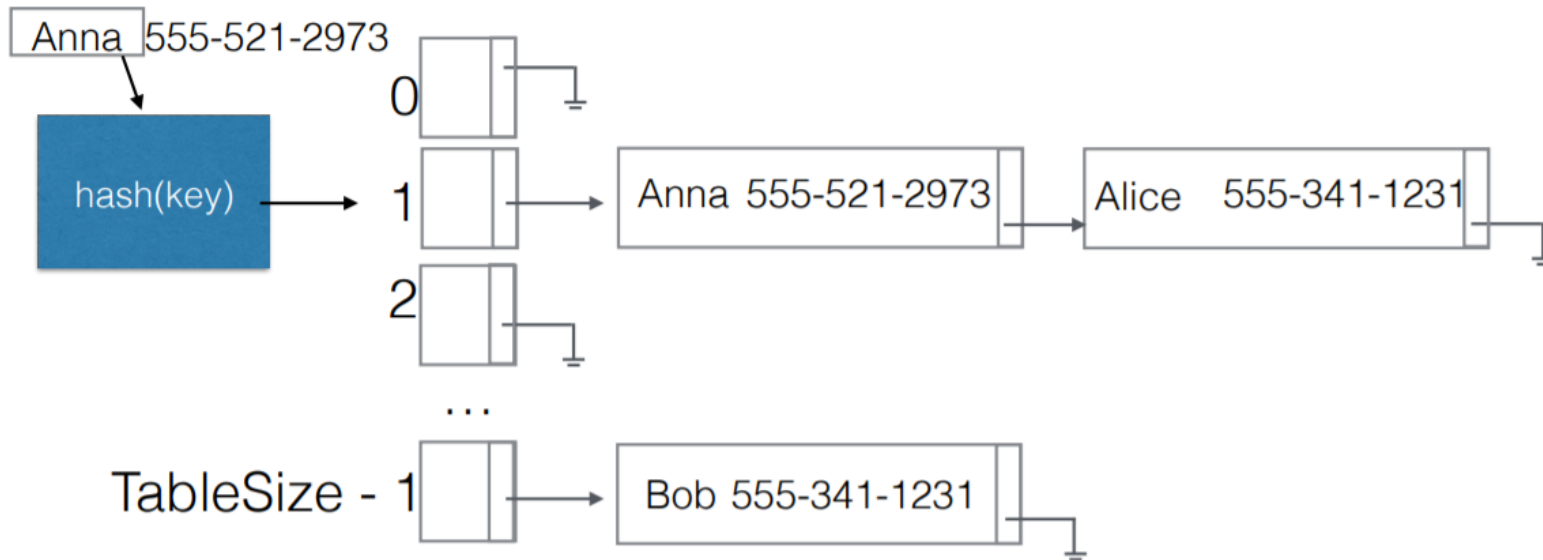
- If we can assume that the possible integer keys  $x$  are distributed evenly, we can use  
 $\text{hash}(x) = x \% \text{TableSize}$
- For strings, can use the ASCII value of the characters
  - E.g.  
 $(\text{str}[n-1] \times 37^{n-1} + \text{str}[n-2] \times 37^{n-2} + \dots + \text{str}[0]) \% \text{TableSize}$

# Choosing a Hash Function

- For a compound object (e.g., a string and an integer), can combine hash functions using
$$\text{hash}(s, x) = (\text{hash}_1(s) \times p_1 + \text{hash}_2(x) \times p_2) \% \text{TableSize}$$
- $p_1, p_2$  prime
  - Don't want to be able to factor the expression
  - $\text{hash}(s, x) = 2 \times (\text{hash}_1(s) \times p_1 + \text{hash}_2(x) \times p_2) \% \text{TableSize}$  means we are skipping the odd cells in the table

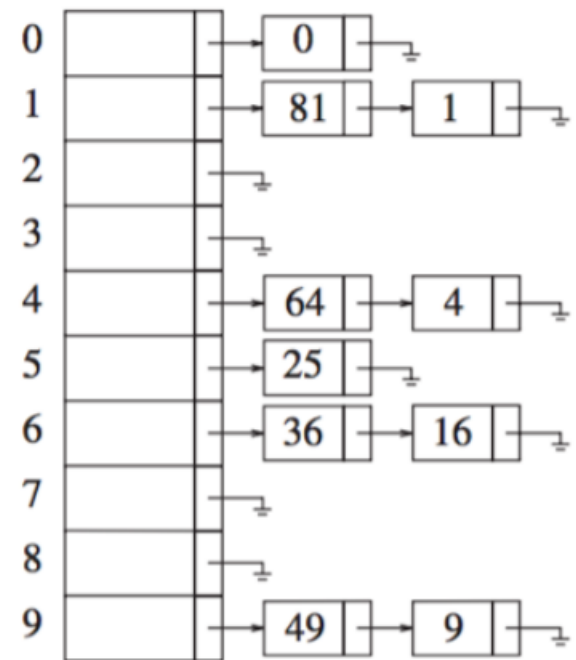
# Collisions: Separate Chaining

- Keep all items whose key hashes to the same value on a linked list
- Can think of each list as a *bucket* defined by the hash value

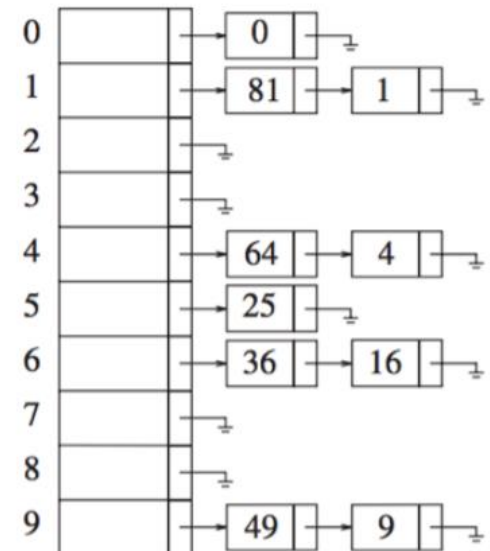


# Runtime for Separate Chaining

- Runtime depends on the number of elements in a list on average
- Load factor  $\lambda = N/TableSize$
- N is the total number of entries

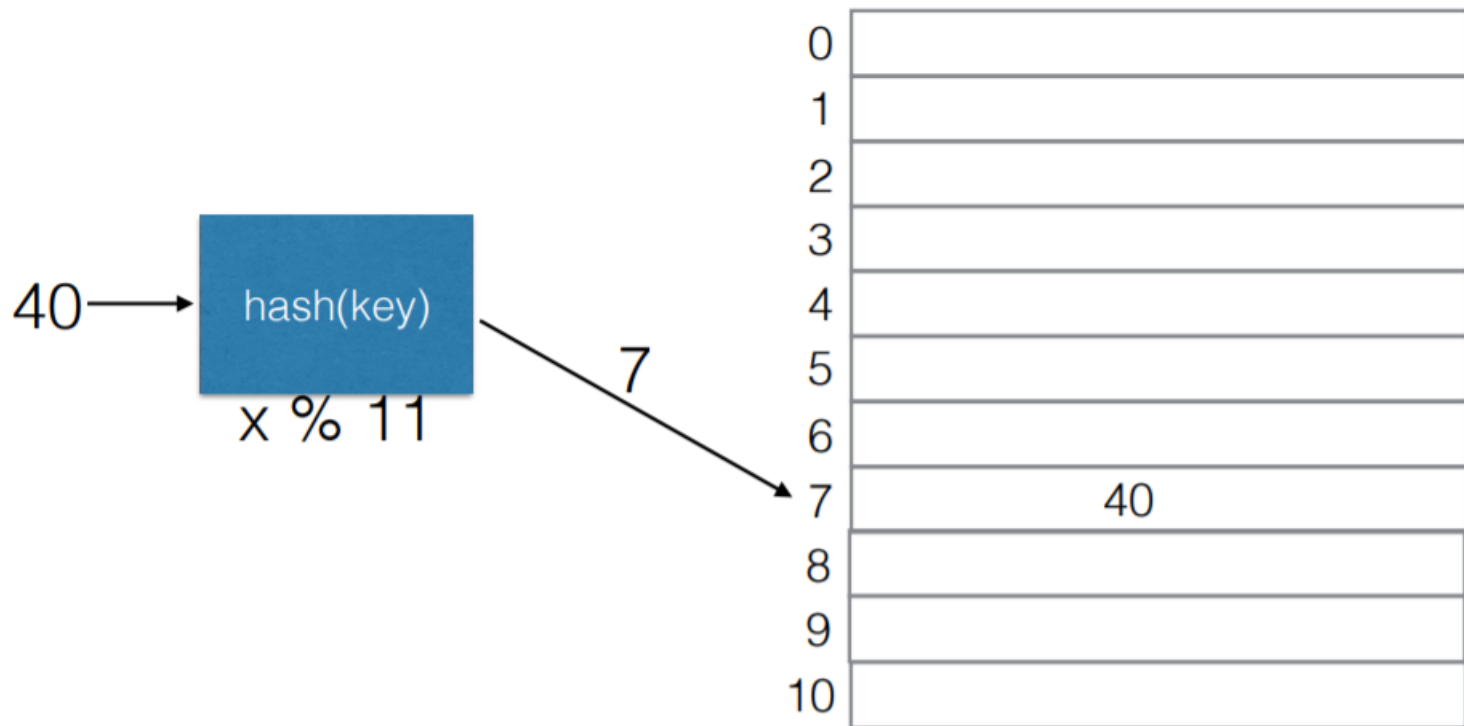


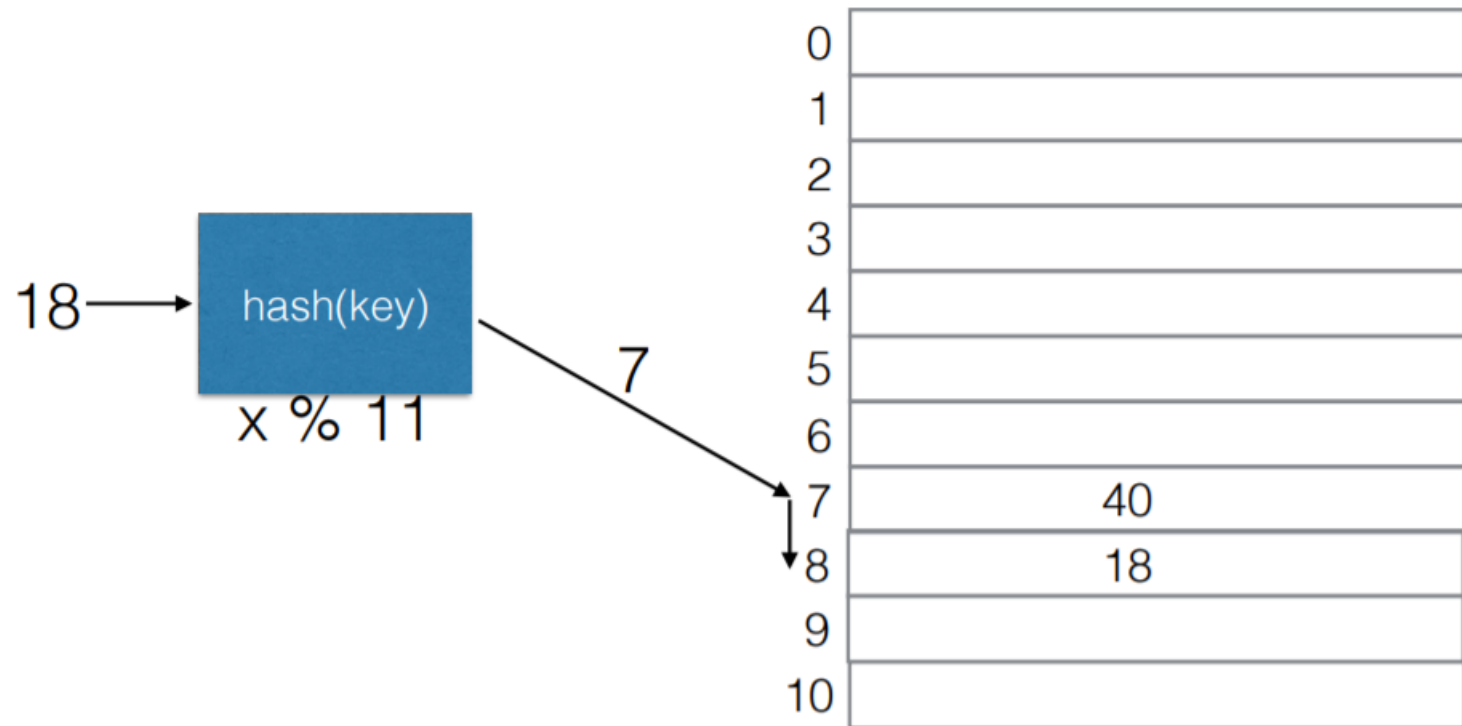
- If lookup fails (key is not in the hash table), need to search all (on average)  $\lambda$  nodes in the list for the hash bucket
- If lookup succeeds (key is in the table), need to search (on average)  $\lambda/2 + 1$  nodes
- Keep  $\lambda \approx 1$ 
  - Increase table size as needed

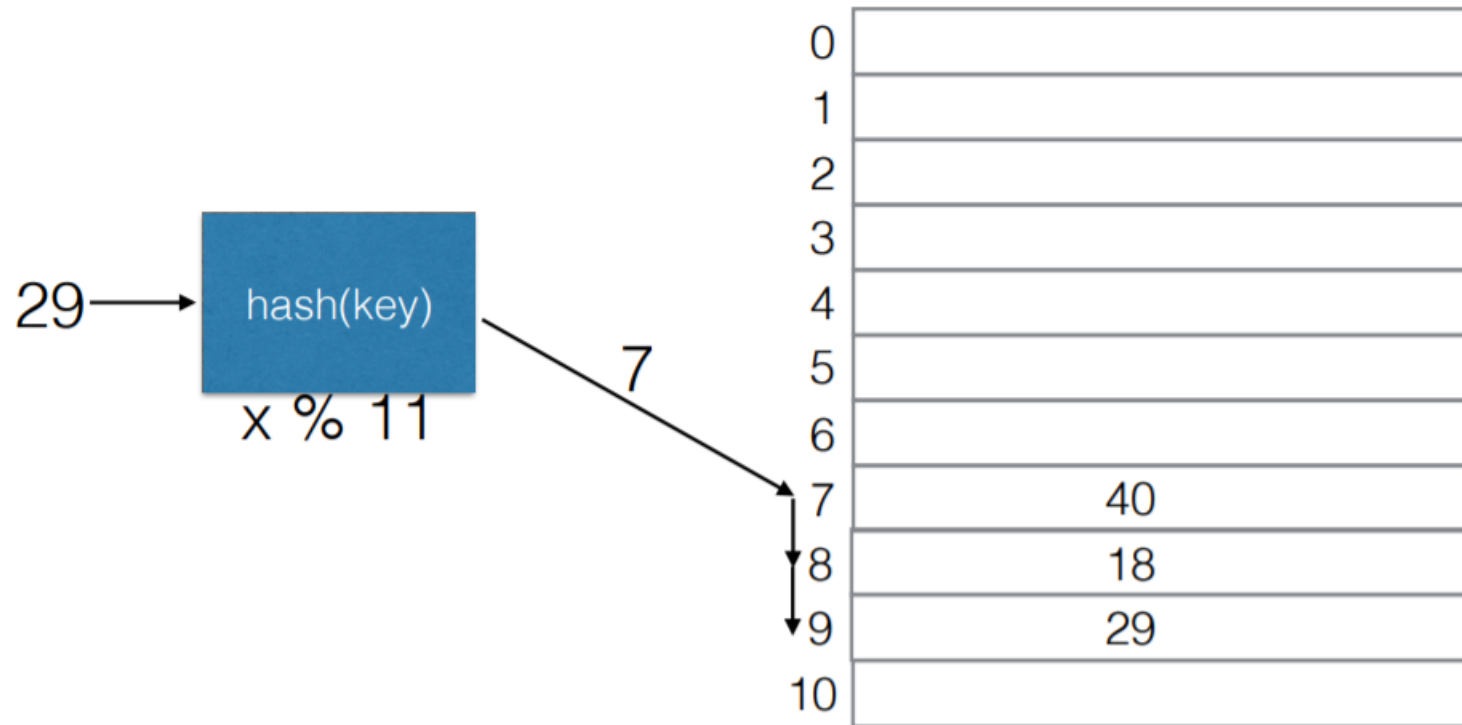


# Probing

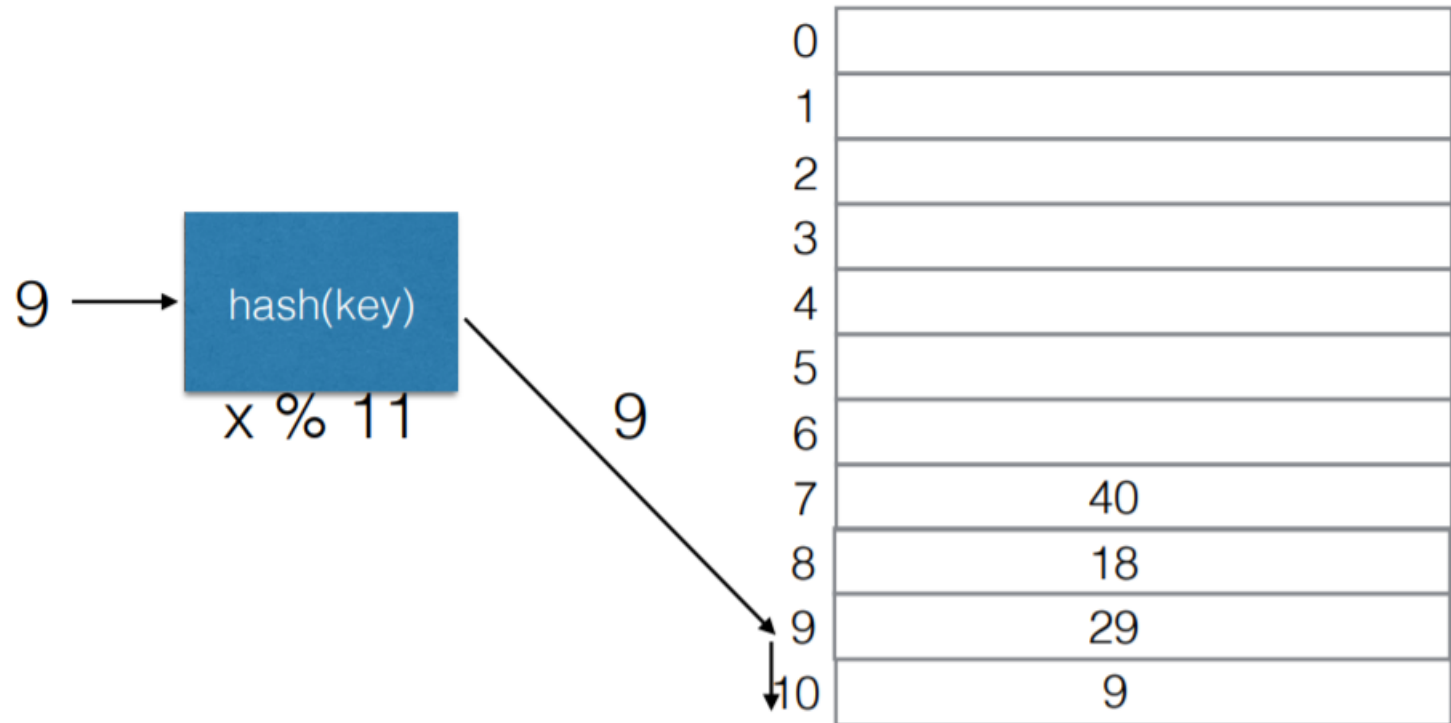
- When a collision occurs, put item in an empty cell of the had table

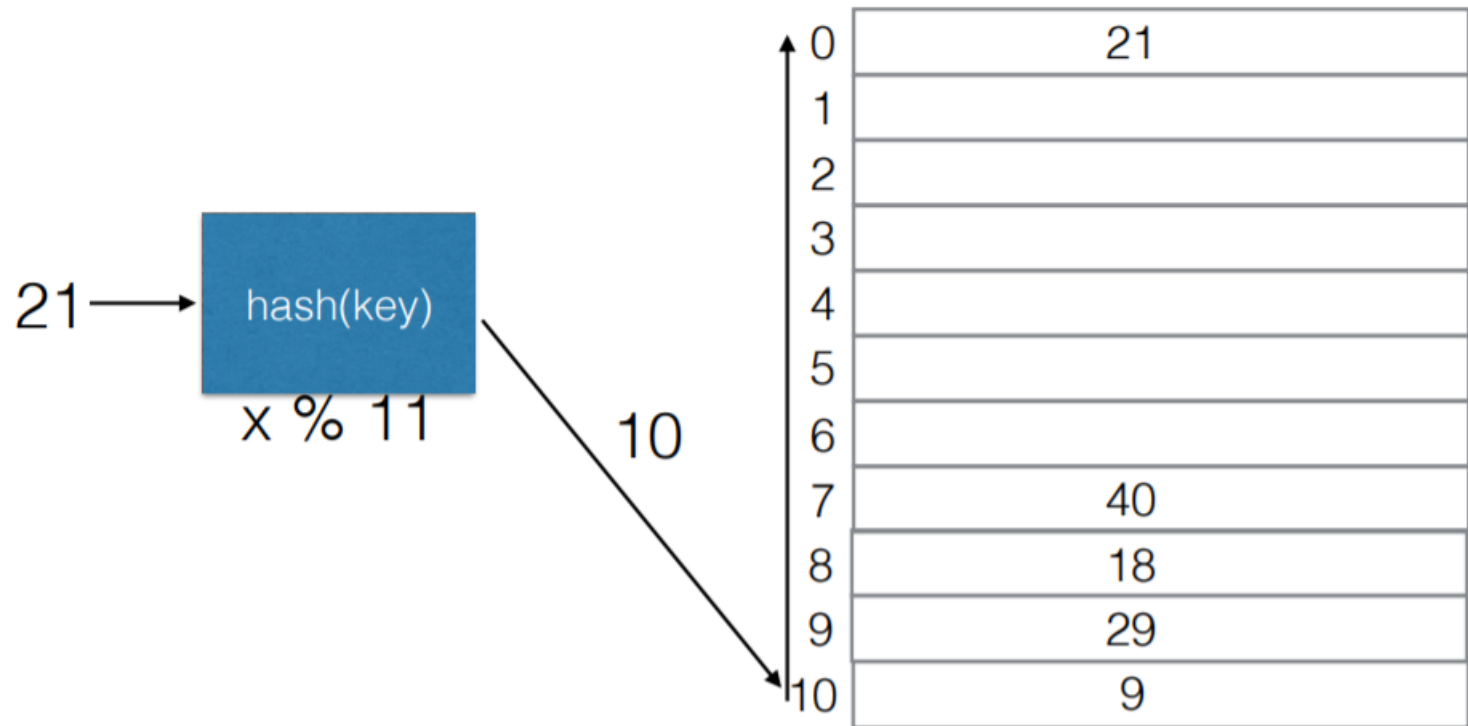




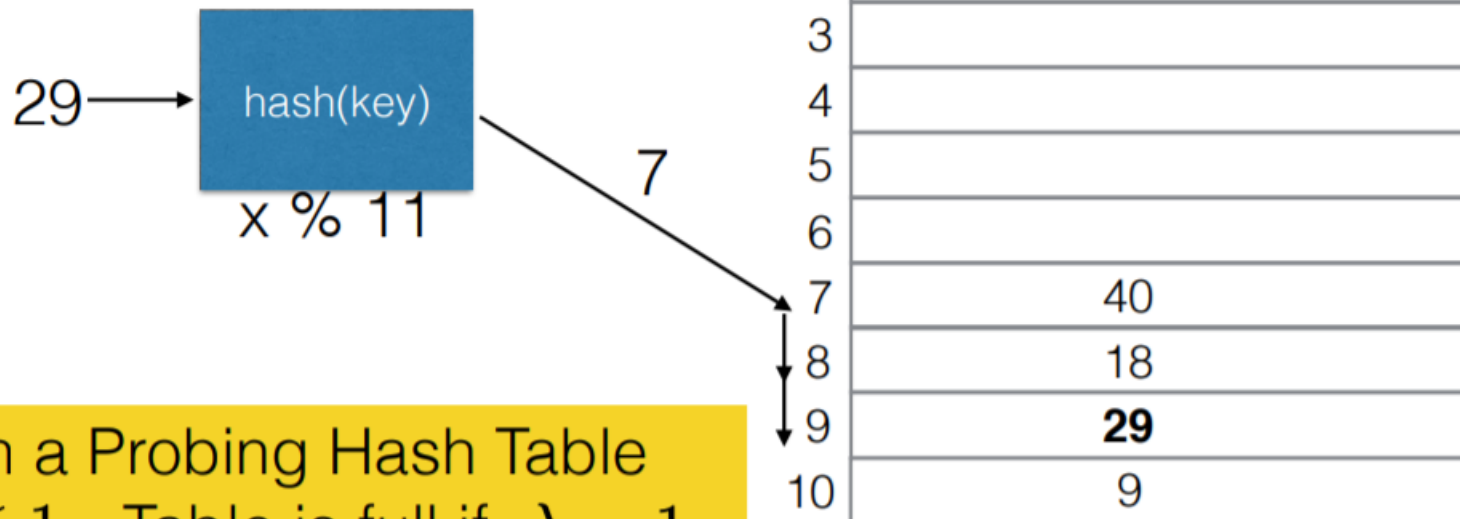








- To look up a key, search the table starting from the cell the key was hashed to



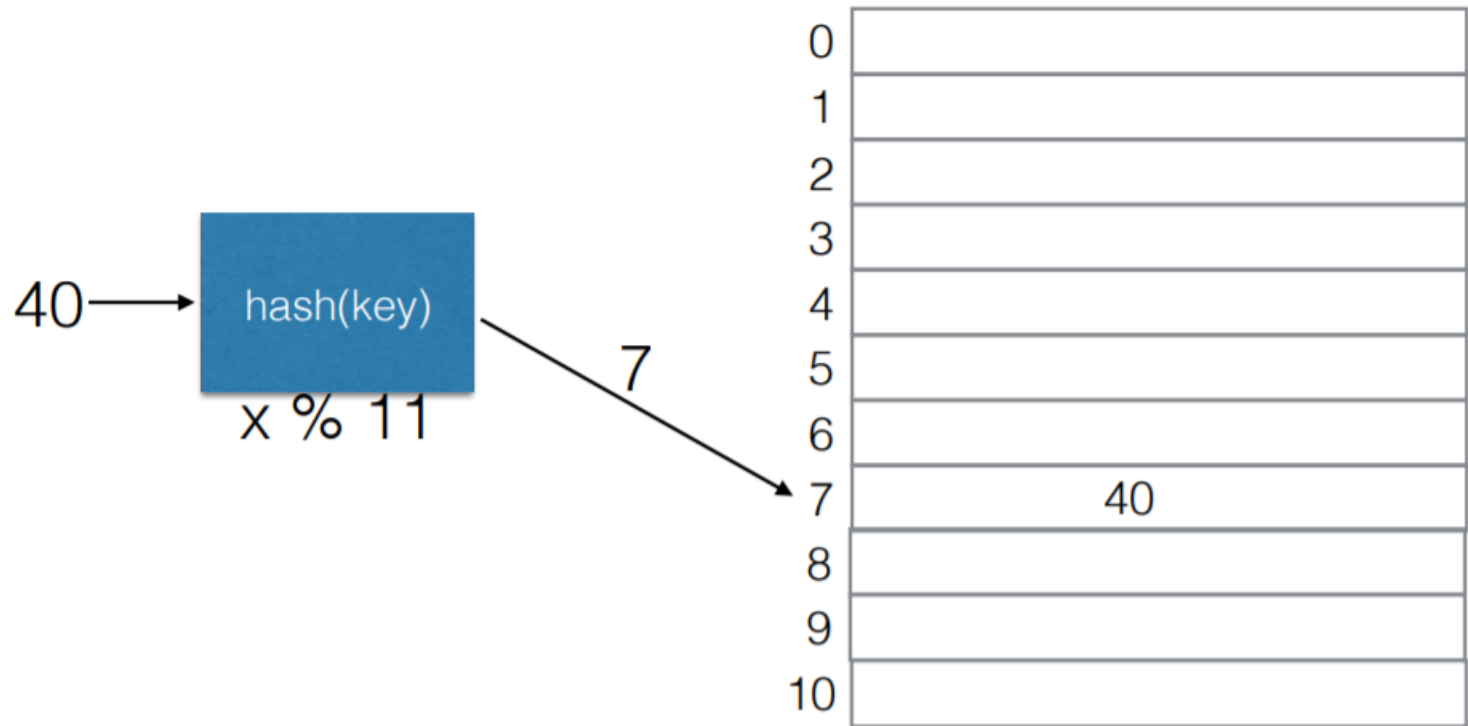
With a Probing Hash Table  
 $\lambda \leq 1$  . Table is full if  $\lambda = 1$  .

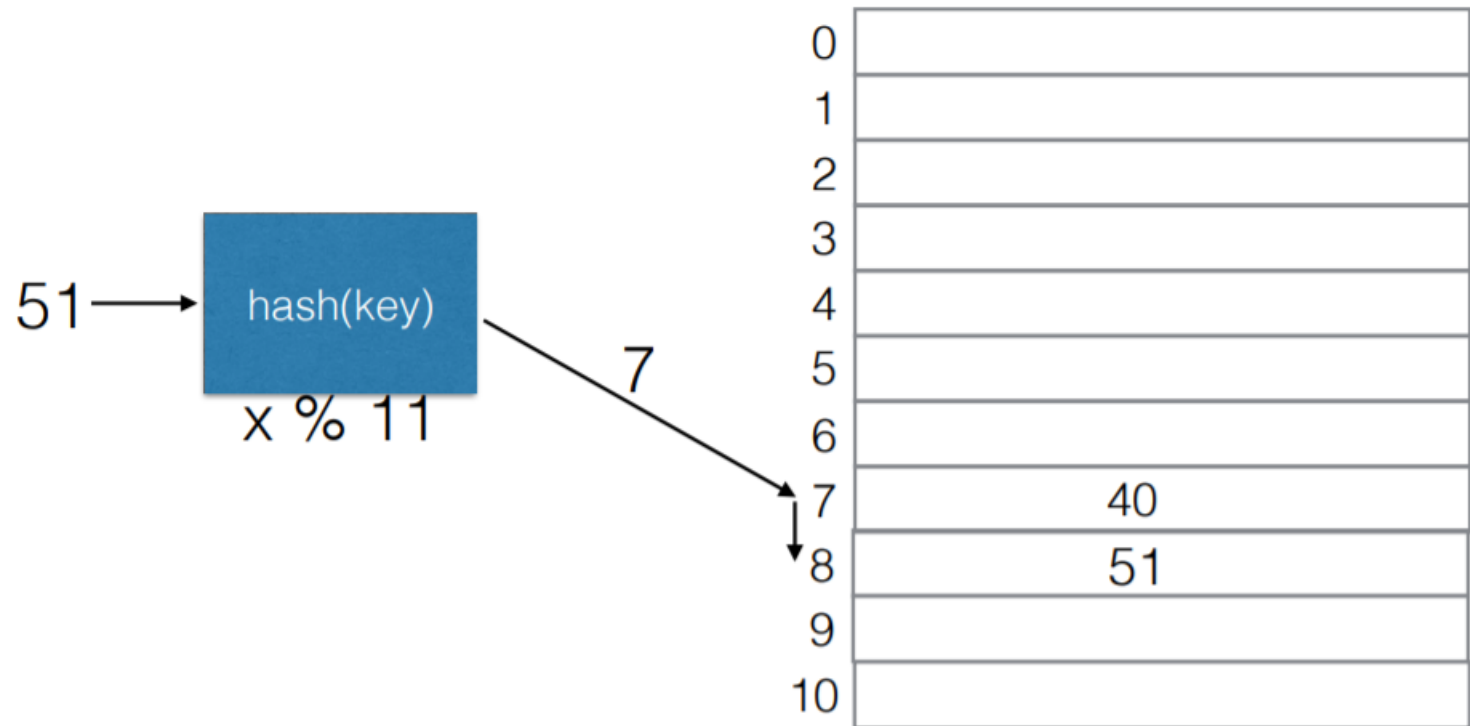
# Probing: Collision Resolution Strategies (1)

- To insert an item, we probe other table cells systematically until an empty cell is found
- To look up a key, we probe until the key is found
- Different strategies to determine the next cell
  - In general, try  $(hash(x) + f(i)) \% TableSize$  after cell  $i$
  - Linear probing: try the next cell,  $f(i) = i$
  - Quadratic probing:  $f(i) = i^2$
  - Double hashing:  $f(i) = i \times hash_2(x)$

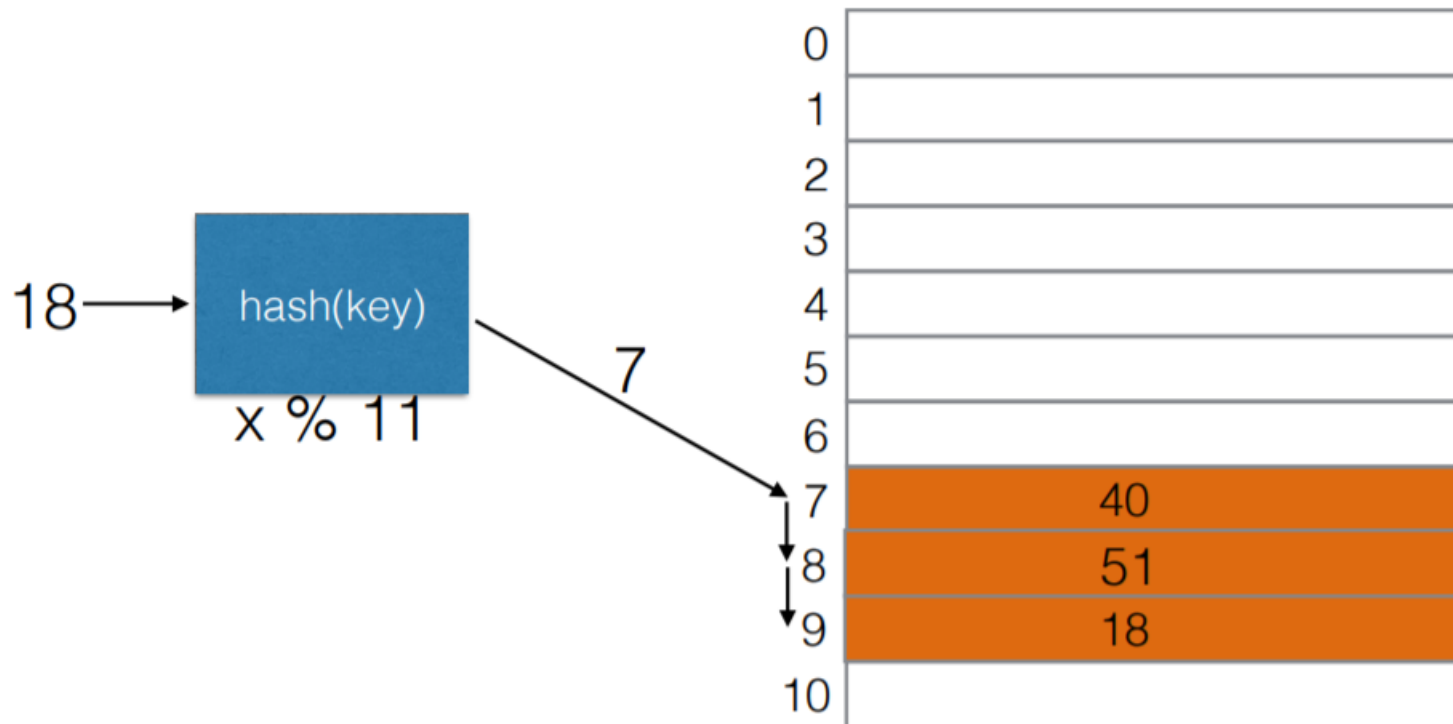
# Linear probing

- If cell is occupied, try the next cell
- Problem: primary clustering
  - Full cells tend to cluster, with no free cells in between
  - Time required to find an empty cell can become very large if the table is almost full ( $\lambda$  is close to 1)

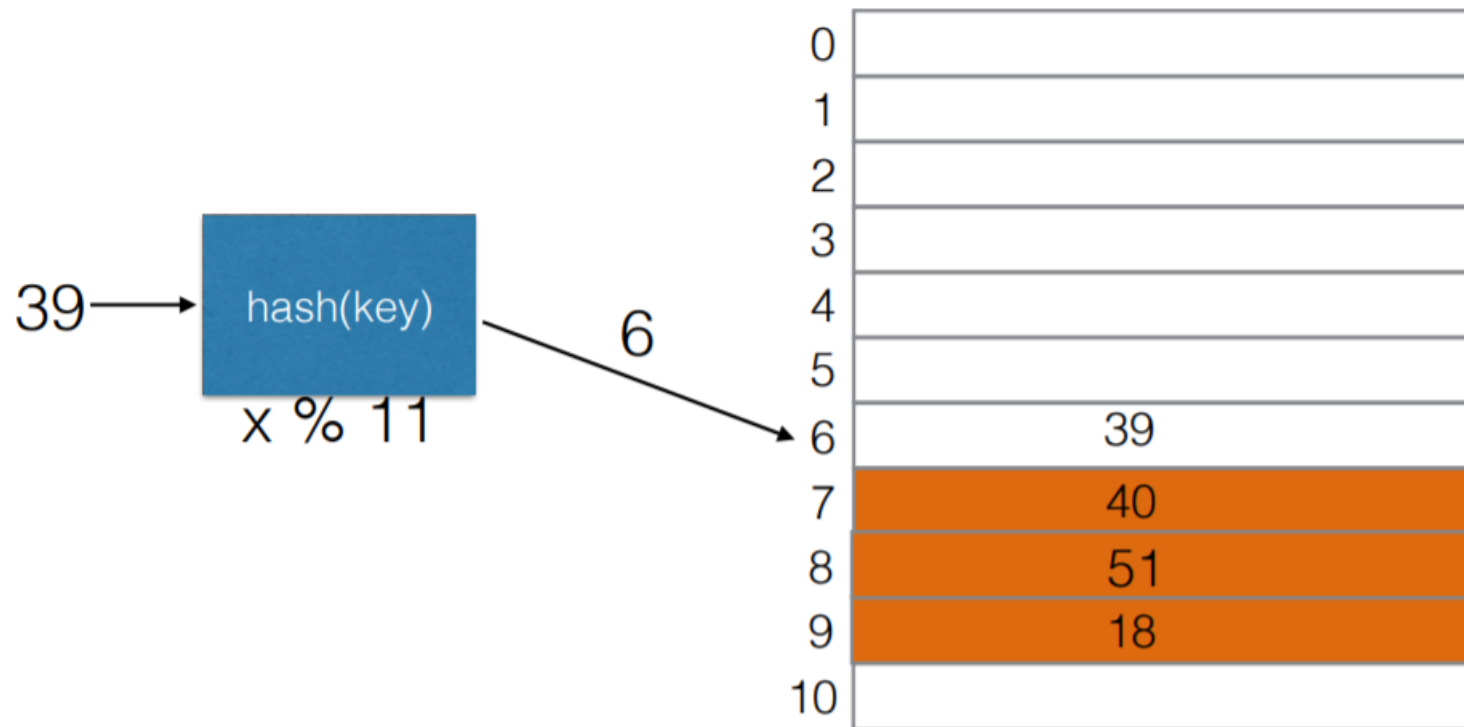


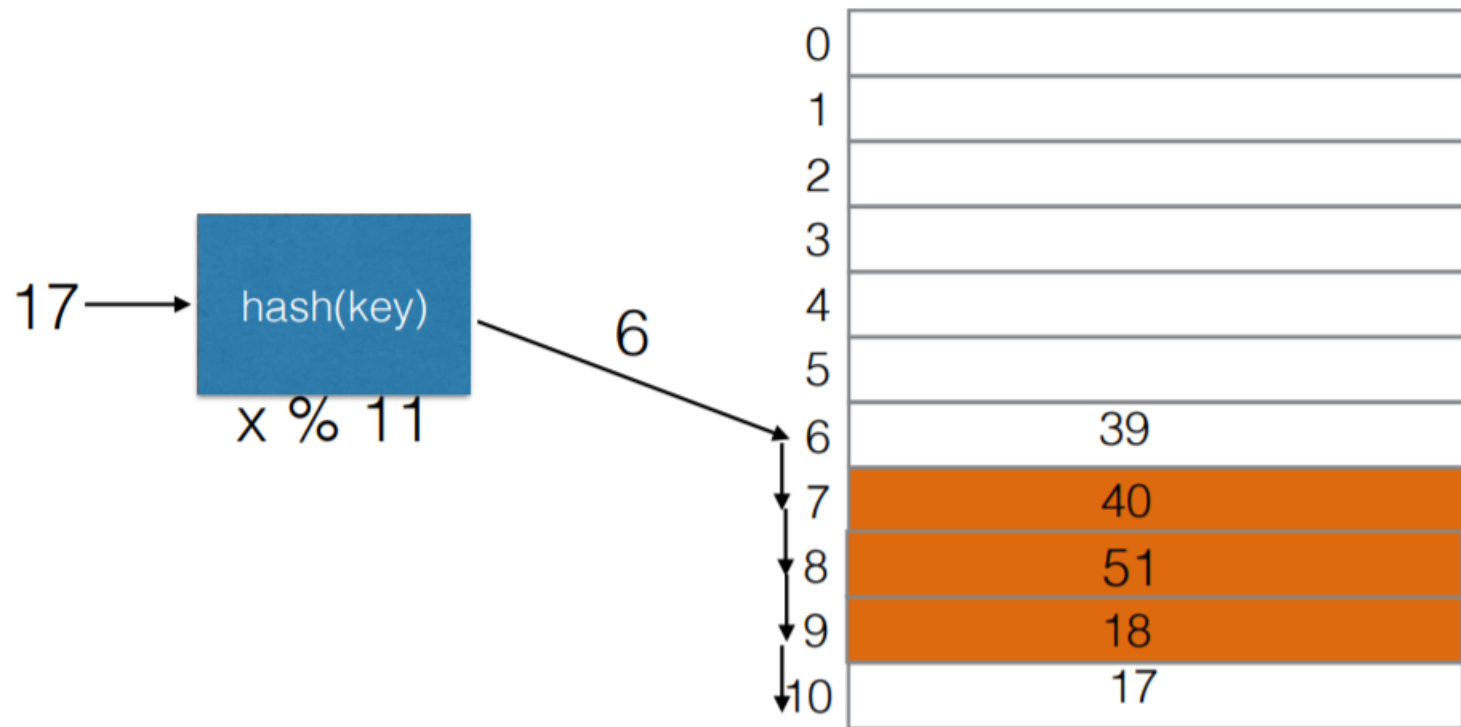


Cells 7-9 are occupied with keys that hash to 7. The entire block is unavailable to keys that hash to  $k < 7$



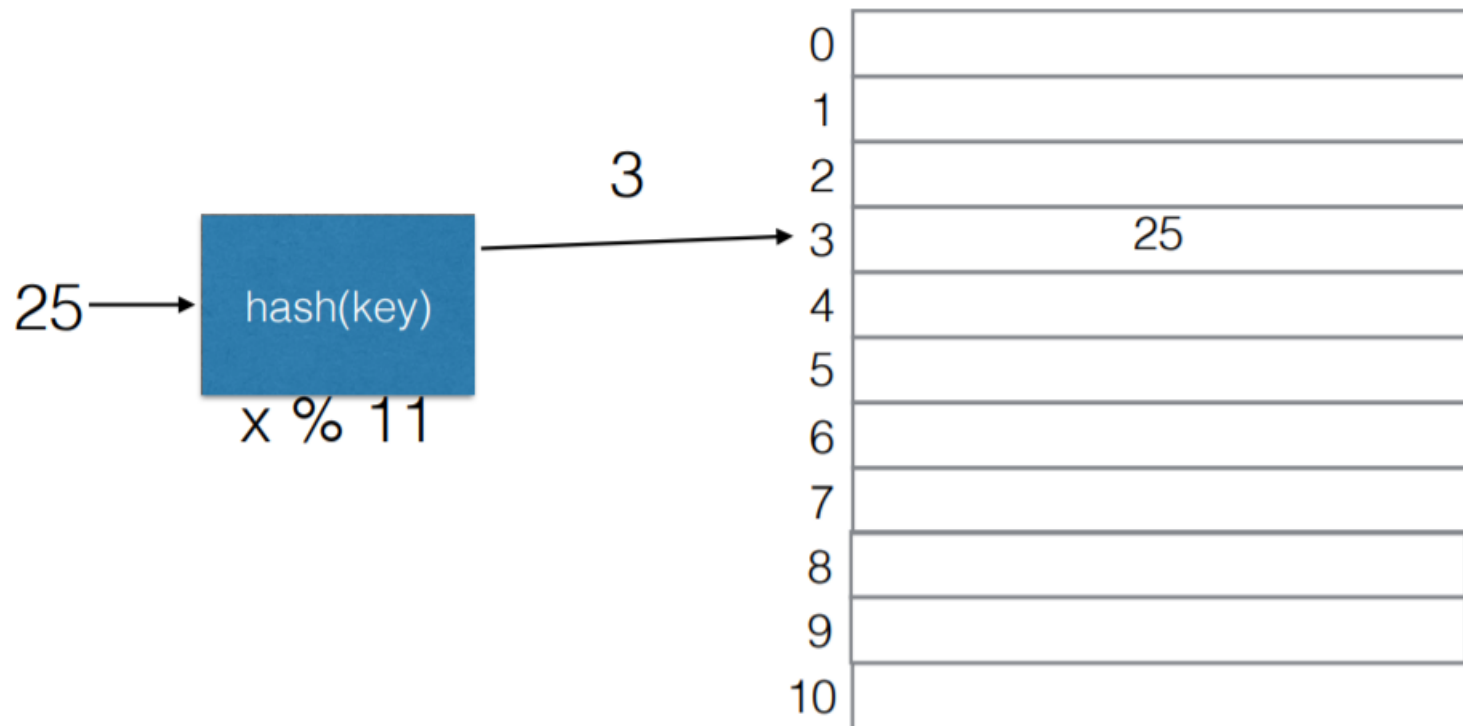


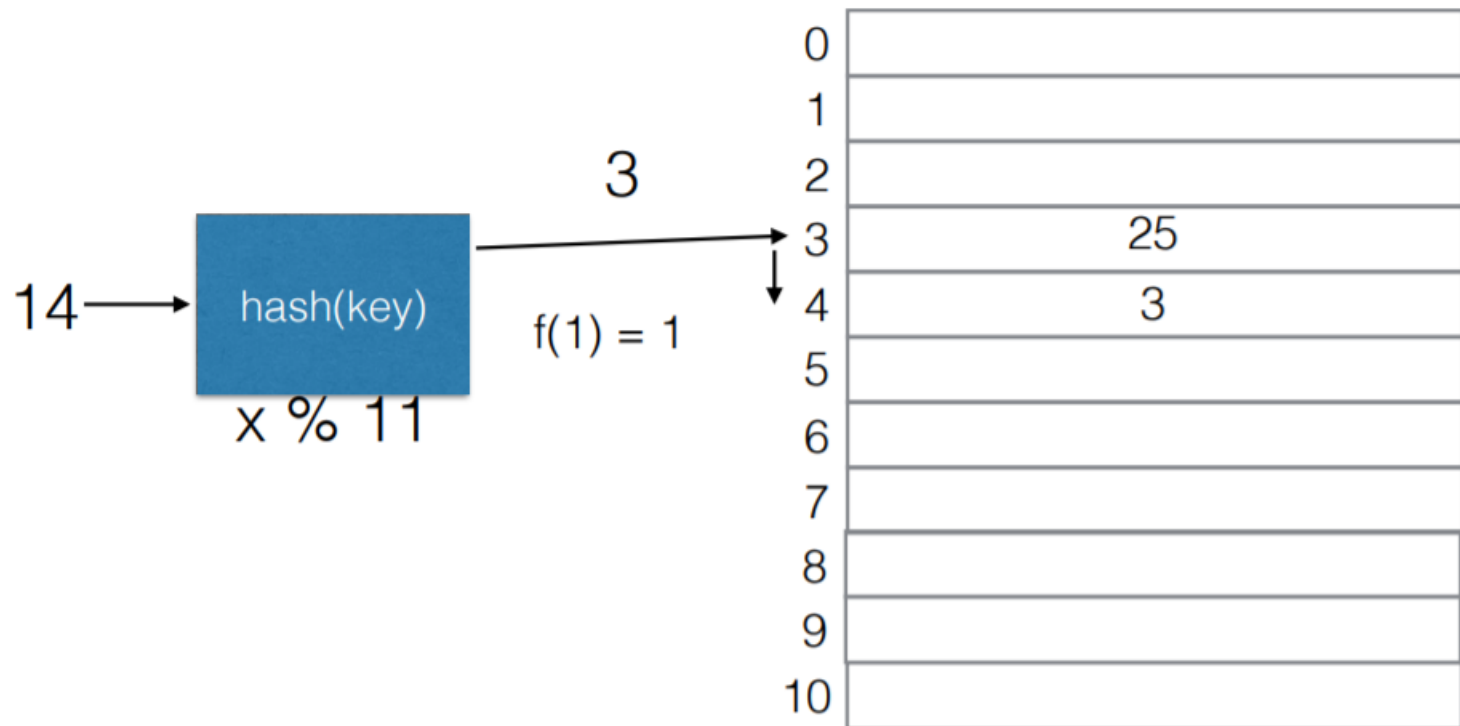


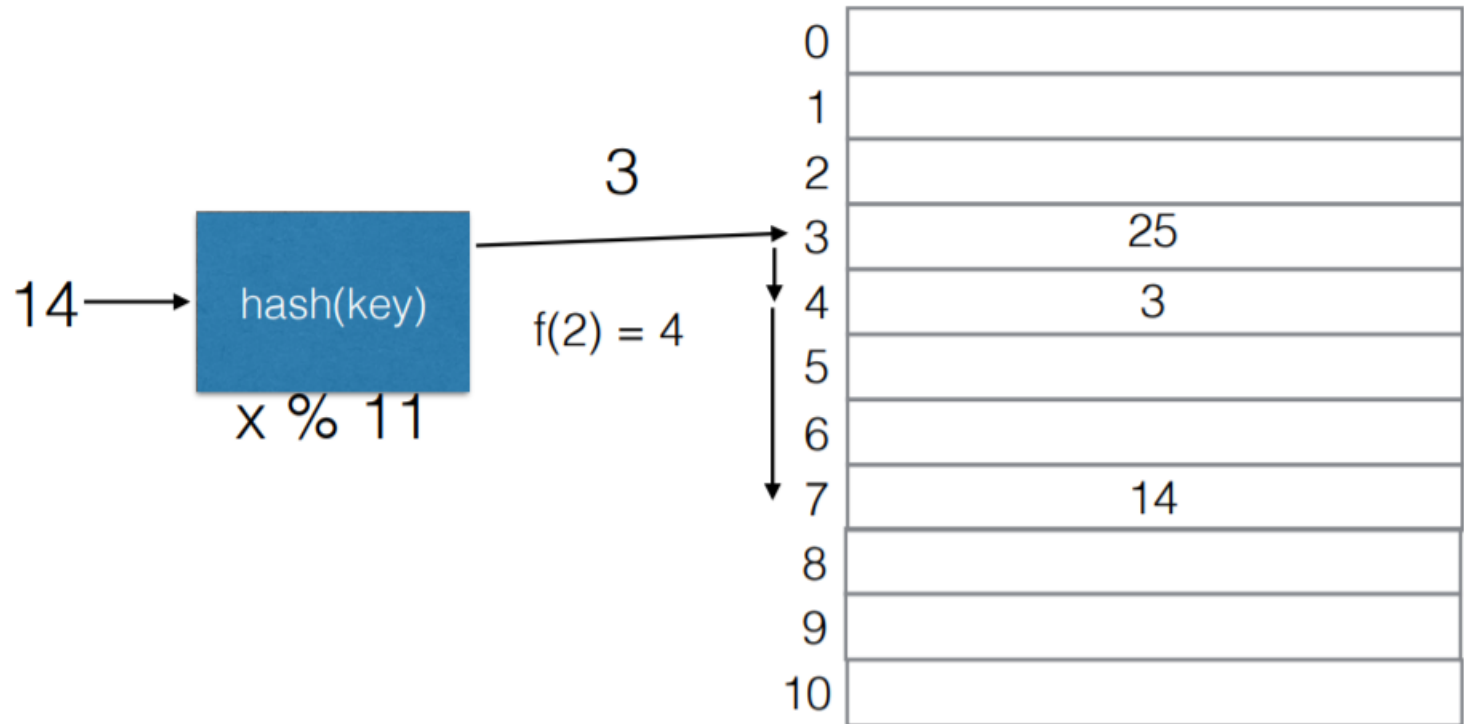


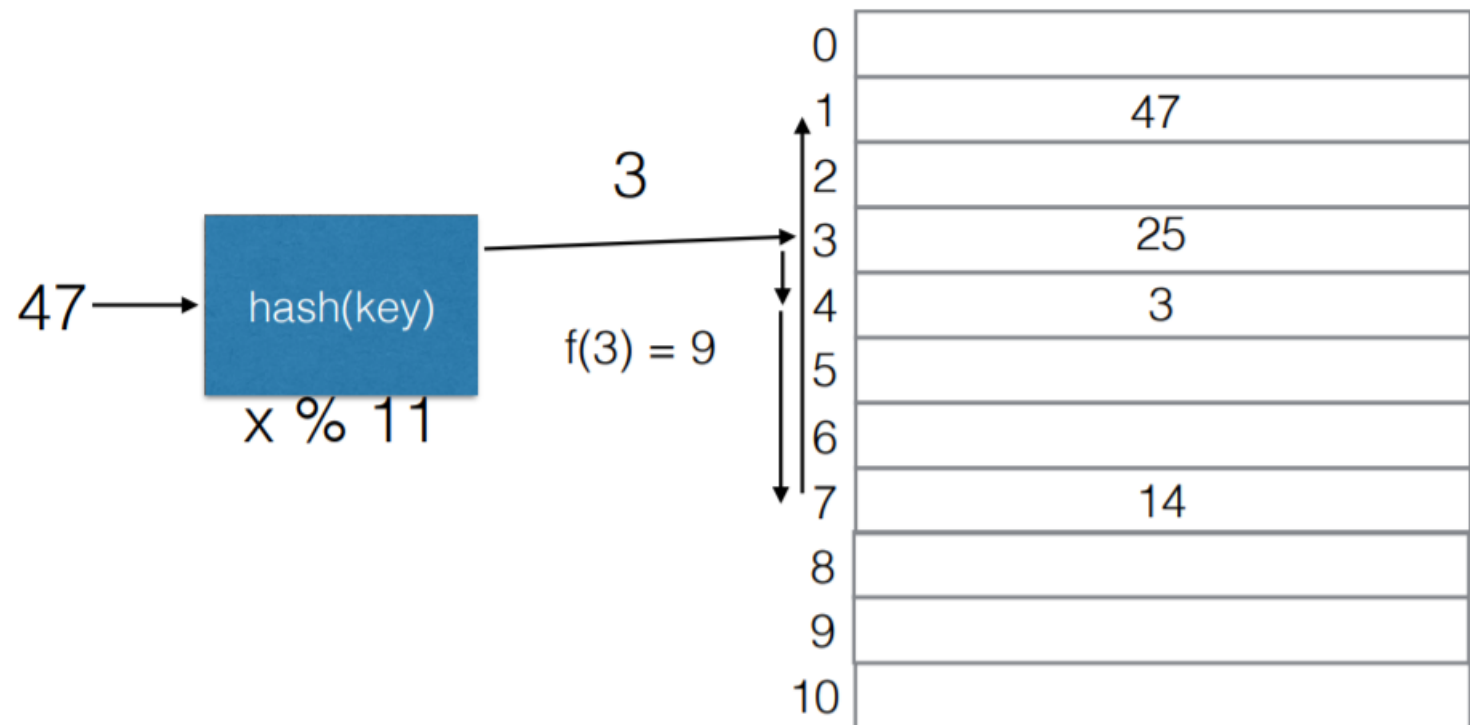
# Quadratic Problem

$$(\text{hash}(x) + f(i)) \% \text{TableSize}, f(i) = i^2$$



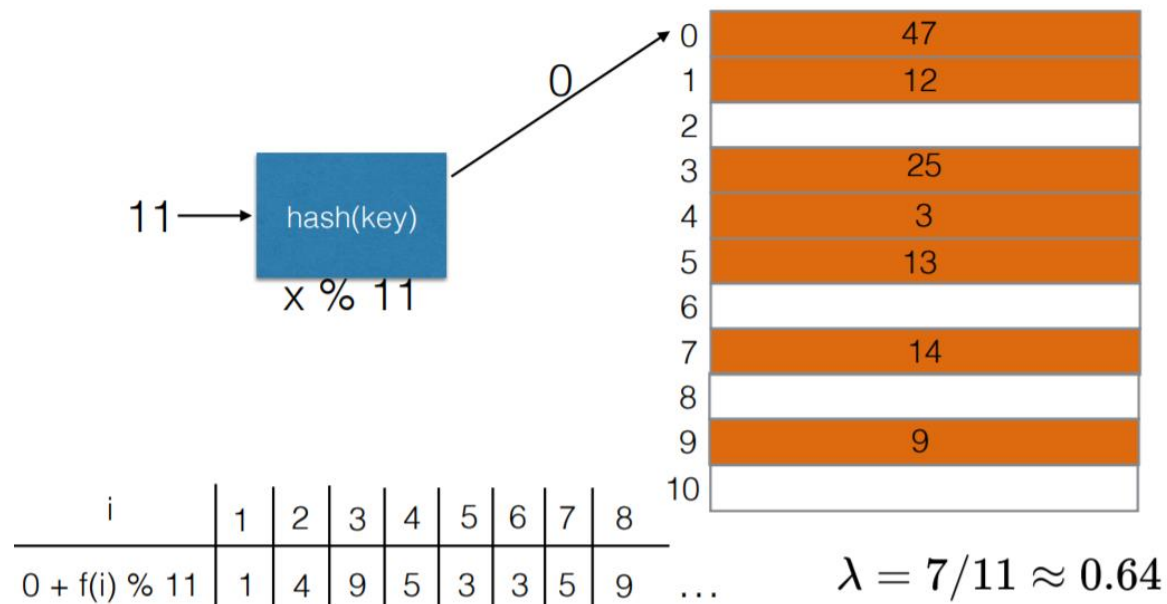






# Quadratic Probing

- Primary clustering is not a problem
- If the table gets too full ( $\lambda > 0.5$ ), it is possible that empty cells become unreachable



# Quadratic Problem Theorem

- If *TableSize* is prime, then the first *TableSize*/2 cells visited by quadratic probing are distinct
  - So it's always possible to find an empty cell if the table is at most half-full

- Proof sketch

Suppose there is a repetition. Then

$$(h + i^2) \% T = (h + j^2) \% T, \text{ so}$$

$$h + i^2 = h + j^2 + kT \text{ for some integer } k, \text{ so}$$

$$kT = (i + j)(i - j)$$

T is prime, so either (i+j) or (i-j) is divisible by T. But  $i < j < M$ , which means that both (i-j) or (i+j) are too small