

Chapter 9

Hash

Data Structures and Algorithms

Luong The Nhan, Tran Giang Son
Faculty of Computer Science and Engineering
University of Technology, VNU-HCM

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

Outcomes

- **L.O.5.1** - Depict the following concepts: hashing table, key, collision, and collision resolution.
- **L.O.5.2** - Describe hashing functions using pseudocode and give examples to show their algorithms.
- **L.O.5.3** - Describe collision resolution methods using pseudocode and give examples to show their algorithms.
- **L.O.5.4** - Implement hashing tables using C/C++.
- **L.O.5.5** - Analyze the complexity and develop experiment (program) to evaluate methods supplied for hashing tables.
- **L.O.1.2** - Analyze algorithms and use Big-O notation to characterize the computational complexity of algorithms composed by using the following control structures: sequence, branching, and iteration (not recursion).



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

Contents

① Basic concepts

② Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

③ Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing



Basic concepts

Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

- Sequential search: $O(n)$
- Binary search: $O(\log_2 n)$

→ Requiring several **key comparisons** before the target is found.



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing



Search complexity:

Size	Binary	Sequential (Average)	Sequential (Worst Case)
16	4	8	16
50	6	25	50
256	8	128	256
1,000	10	500	1,000
10,000	14	5,000	10,000
100,000	17	50,000	100,000
1,000,000	20	500,000	1,000,000

Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing



Is there a search algorithm
whose complexity is $O(1)$?

Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing



Is there a search algorithm
whose complexity is $O(1)$?

YES

Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing



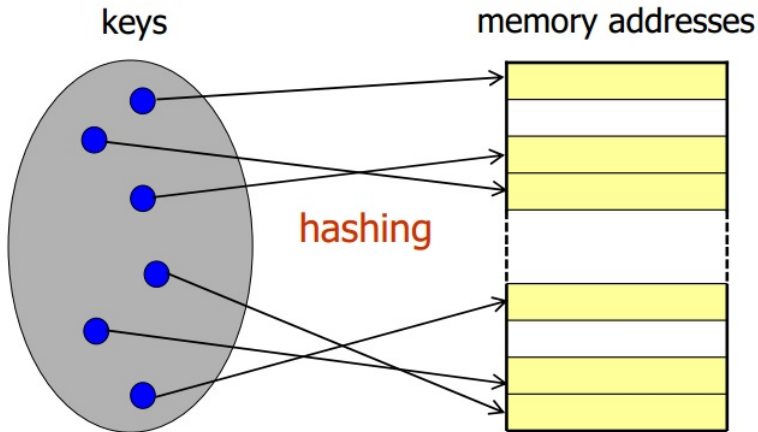
Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

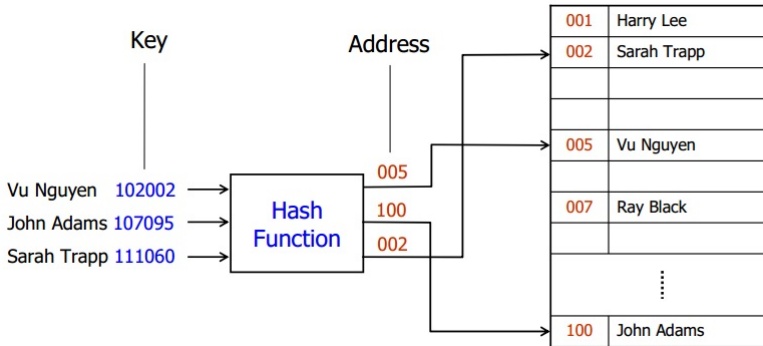
Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing



Hình: Each key has only one address

Basic concepts



Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

Basic concepts

- **Home address**: address produced by a hash function.
- **Prime area**: memory that contains all the home addresses.

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

Basic concepts

- **Home address**: address produced by a hash function.
- **Prime area**: memory that contains all the home addresses.
- **Synonyms**: a set of keys that hash to the same location.
- **Collision**: the location of the data to be inserted is already occupied by the synonym data.



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

Basic concepts

- **Home address**: address produced by a hash function.
- **Prime area**: memory that contains all the home addresses.
- **Synonyms**: a set of keys that hash to the same location.
- **Collision**: the location of the data to be inserted is already occupied by the synonym data.
- **Ideal hashing**:
 - No location collision
 - Compact address space



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

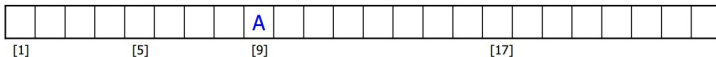
Basic concepts

Insert A, B, C

hash(A) = 9

hash(B) = 9

hash(C) = 17



Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

Basic concepts

Insert A, B, C

hash(A) = 9

hash(B) = 9

hash(C) = 17

B and A
collide at 9



Collision Resolution



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

Basic concepts

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

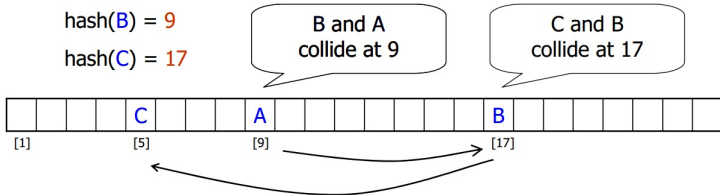
- Open addressing
- Linked list resolution
- Bucket hashing

Insert A, B, C

hash(A) = 9

hash(B) = 9

hash(C) = 17



Collision Resolution

Basic concepts

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

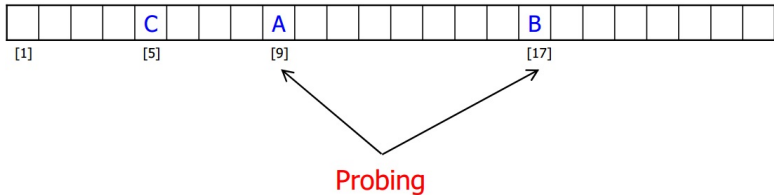
- Open addressing
- Linked list resolution
- Bucket hashing

Search for B

$\text{hash}(A) = 9$

$\text{hash}(B) = 9$

$\text{hash}(C) = 17$





Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

Hash functions

Hash functions

- Direct hashing
- Modulo division
- Digit extraction
- Mid-square
- Folding
- Rotation
- Pseudo-random

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

Direct Hashing

The address is the key itself:

$$\text{hash}(\text{Key}) = \text{Key}$$

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

Direct Hashing

- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

- **Advantage:** there is no collision.
- **Disadvantage:** the address space (storage size) is as large as the key space.



$$Address = Key \bmod listSize$$

- Fewer collisions if *listSize* is a prime number.
- Example:
Numbering system to handle 1,000,000 employees
Data space to store up to 300 employees
 $hash(121267) = 121267 \bmod 307 = 2$



Digit extraction

Address = selected digits from Key

Example:

379452 → 394

121267 → 112

378845 → 388

160252 → 102

045128 → 051

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

Direct Hashing

Modulo division

Digit extraction

Mid-square

Mid-square

Folding

Rotation

Pseudo-random

Collision resolution

Open addressing

Linked list resolution

Bucket hashing

Mid-square

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction

Mid-square

- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

Address = middle digits of Key^2

Example:

$$9452 * 9452 = 89340304 \rightarrow 3403$$

Mid-square

- **Disadvantage:** the size of the Key^2 is too large.
- **Variations:** use only a portion of the key.

Example:

$$379452: 379 * 379 = 143641 \rightarrow 364$$

$$121267: 121 * 121 = 014641 \rightarrow 464$$

$$045128: 045 * 045 = 002025 \rightarrow 202$$



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square

Mid-square

Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

Folding

The key is divided into parts whose size matches the address size.

Example:

Key = 123|456|789

fold shift

$123 + 456 + 789 = 1368$

→ 368

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square

Folding

- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

Folding

The key is divided into parts whose size matches the address size.

Example:

Key = 123|456|789

fold shift

$$123 + 456 + 789 = 1368$$

→ 368

fold boundary

$$321 + 456 + 987 = 1764$$

→ 764



- Hashing keys that are identical except for the last character may create synonyms.
- The key is rotated before hashing.

original key	rotated key
600101	160010
600102	260010
600103	360010
600104	460010
600105	560010



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding

Rotation

Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

- Used in combination with fold shift.

original key	rotated key
600101 → 62	160010 → 26
600102 → 63	260010 → 36
600103 → 64	360010 → 46
600104 → 65	460010 → 56
600105 → 66	560010 → 66

Spreading the data more evenly across the address space.



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding

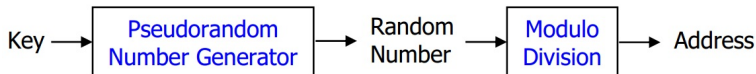
Rotation

Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

Pseudo-random



$$y = ax + c$$

For maximum efficiency, a and c should be prime numbers.



Basic concepts

Hash functions

Direct Hashing

Modulo division

Digit extraction

Mid-square

Mid-square

Folding

Rotation

Pseudo-random

Collision resolution

Open addressing

Linked list resolution

Bucket hashing

Pseudo-random

Example:

Key = 121267

$a = 17$

$c = 7$

listSize = 307

Address = $((17 * 121267 + 7) \bmod 307$

$= (2061539 + 7) \bmod 307$

$= 2061546 \bmod 307$

$= 41$

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation

Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

Collision resolution

- Except for the direct hashing, none of the others are **one-to-one mapping**
→ Requiring collision resolution methods
- Each collision resolution method can be used **independently** with each hash function



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

- A rule of thumb: a hashed list should not be allowed to become more than 75% full.

Load factor:

$$\alpha = (k/n) \times 100$$

n = list size

k = number of filled elements



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

- As data are added and collisions are resolved, hashing tends to cause data to group within the list.
→ **Clustering**: data are unevenly distributed across the list.
- High degree of clustering increases the **number of probes** to locate an element.
→ **Minimize** clustering.



Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Open addressing
Linked list resolution
Bucket hashing



Basic concepts

Hash functions

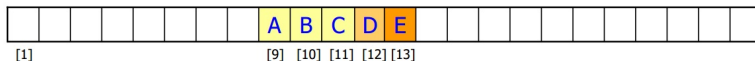
Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

- Primary clustering: data become clustered around a home address.

Insert A_9 , B_9 , C_9 , D_{11} , E_{12}





Basic concepts

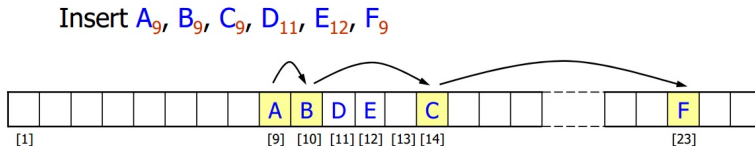
Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

- **Secondary clustering:** data become grouped along a collision path throughout a list.



- Open addressing
- Linked list resolution
- Bucket hashing



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

Open addressing

- Linked list resolution
- Bucket hashing

When a collision occurs, an **unoccupied element** is searched for placing the new element in.



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

Open addressing

- Linked list resolution
- Bucket hashing

Hash function:

$$h : U \rightarrow \{0, 1, 2, \dots, m - 1\}$$

set of keys

addresses

Open addressing

Hash and probe function:

$$hp : U \times \{0, 1, 2, \dots, m-1\} \rightarrow \{0, 1, 2, \dots, m-1\}$$

set of keys

probe numbers

addresses

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

Open addressing

- Linked list resolution
- Bucket hashing

Open Addressing

Algorithm hashInsert(ref T <array>, val k <key>)
Inserts key k into table T.

```
i = 0
while i < m do
    j = hp(k, i)
    if T[j] = nil then
        T[j] = k
        return j
    else
        i = i + 1
    end
end
return error: "hash table overflow"
End hashInsert
```



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

Open addressing

- Linked list resolution
- Bucket hashing

Open Addressing

Algorithm hashSearch(val T <array>, val k <key>)
Searches for key k in table T.

```
i = 0
while i < m do
    j = hp(k, i)
    if T[j] = k then
        | return j
    else if T[j] = nil then
        | return nil
    else
        | i = i + 1
    end
end
return nil
End hashSearch
```



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

Open addressing

- Linked list resolution
- Bucket hashing



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing

Linked list resolution
Bucket hashing

There are different methods:

- Linear probing
- Quadratic probing
- Double hashing
- Key offset

Linear Probing

- When a home address is occupied, go to the **next address** (the current address + 1):
$$hp(k, i) = (h(k) + i) \bmod m$$



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing

Linked list resolution
Bucket hashing

Hash

**Luong The Nhan,
Tran Giang Son**



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

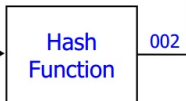
Collision resolution

- Open addressing
- Linked list resolution
- Bucket hashing

-
- Harry Eagle 166702 → Hash Function → 002
- | | | |
|-----|----------------|----------|
| 001 | Mary Dodd | (379452) |
| 002 | Sarah Trapp | (070918) |
| 003 | Bryan Devaux | (121267) |
| | | |
| | | |
| | | |
| | | |
| 008 | John Carver | (378845) |
| | | |
| ⋮ | | |
| 306 | Tuan Ngo | (160252) |
| 307 | Shouli Feldman | (045128) |

Linear Probing

Harry Eagle 166702



001	Mary Dodd	(379452)
002	Sarah Trapp	(070918)
003	Bryan Devaux	(121267)
004	Harry Eagle	(166702)
008	John Carver	(378845)
⋮		
306	Tuan Ngo	(160252)
307	Shouli Feldman	(045128)

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

Open addressing

- Linked list resolution
- Bucket hashing

- Advantages:
 - quite simple to implement
 - data tend to remain near their home address (significant for disk addresses)
- Disadvantages:
 - produces primary clustering



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing

Linked list resolution
Bucket hashing



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing

Linked list resolution
Bucket hashing

- The address increment is the **collision probe number** squared:
$$hp(k, i) = (h(k) + i^2) \bmod m$$

- Advantages:
 - works much better than linear probing
 - Disadvantages:
 - time required to square numbers
 - produces secondary clustering
- $$h(k_1) = h(k_2) \rightarrow hp(k_1, i) = hp(k_2, i)$$





Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing

Linked list resolution
Bucket hashing

- Using **two** hash functions:
$$hp(k, i) = (h_1(k) + ih_2(k)) \bmod m$$

- The new address is a function of the collision address and the key.

$$\begin{aligned} offset &= [key / listSize] \\ newAddress &= \\ (collisionAddress + offset) \bmod listSize \end{aligned}$$



- The new address is a function of the collision address and the key.

$$\begin{aligned} offset &= [key / listSize] \\ newAddress &= \\ & (collisionAddress + offset) \bmod listSize \end{aligned}$$

$$hp(k, i) = (hp(k, i - 1) + [k/m]) \bmod m$$





Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

Open addressing

- Linked list resolution
- Bucket hashing

Hash and probe function:

$$hp : U \times \{0, 1, 2, \dots, m-1\} \rightarrow \{0, 1, 2, \dots, m-1\}$$

set of keys probe numbers addresses

$\{hp(k, 0), hp(k, 1), \dots, hp(k, m-1)\}$ is a permutation of $\{0, 1, \dots, m-1\}$



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution
Bucket hashing

- Major disadvantage of Open Addressing:
each collision resolution increases the
probability for future collisions.
→ use **linked lists** to store synonyms

Linked list resolution

001	Mary Dodd	(379452)	×
002	Sarah Trapp	(070918)	→
003	Bryan Devaux	(121267)	×
			×
			×
			×
			×
008	John Carver	(378845)	×
			×
⋮			
306	Tuan Ngo	(160252)	×
307	Shouli Feldman	(045128)	×



overflow area

prime area



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

Open addressing

Linked list resolution

Bucket hashing

- Hashing data to **buckets** that can hold multiple pieces of data.
- Each bucket has an address and **collisions** are **postponed** until the bucket is full.



Basic concepts

Hash functions

Direct Hashing
Modulo division
Digit extraction
Mid-square
Mid-square
Folding
Rotation
Pseudo-random

Collision resolution

Open addressing
Linked list resolution

Bucket hashing

Bucket hashing

001	Mary Dodd (379452)
002	Sarah Trapp (070918)
	Harry Eagle (166702)
	Ann Georgis (367173)
003	Bryan Devaux (121267)
	Chris Walljasper(572556)
⋮	
307	Shouli Feldman (045128)



linear probing

Hash

Luong The Nhan,
Tran Giang Son



Basic concepts

Hash functions

- Direct Hashing
- Modulo division
- Digit extraction
- Mid-square
- Mid-square
- Folding
- Rotation
- Pseudo-random

Collision resolution

- Open addressing
- Linked list resolution

Bucket hashing