

# Universal Hashing

COMPSI 753: Algorithms for Massive Data

Ninh Pham

University of Auckland

Parts of this material are modifications of the lecture slides of Kevin Wayne

(<https://www.cs.princeton.edu/~wayne/kleinberg-tardos/>)

Designed for the textbook **Algorithm Design**

by Jon Kleinberg and Eva Tardos.

Auckland, Aug 3, 2020

# Dictionary data type

- Dictionary problem:

- Given a universe  $\mathbf{U}$  of possible elements, maintain a subset  $\mathbf{S} \subseteq \mathbf{U}$  supporting fast searching, inserting, deleting in  $\mathbf{S}$  (in **constant** time).

- Main dictionary operations:

- **Insert( $\mathbf{x}$ )**: insert element  $\mathbf{x} \in \mathbf{U}$  to  $\mathbf{S}$ .
- **Delete( $\mathbf{x}$ )**: remove element  $\mathbf{x}$  from  $\mathbf{S}$  if  $\mathbf{x} \in \mathbf{S}$ .
- **Search( $\mathbf{x}$ )**: return true if  $\mathbf{x} \in \mathbf{S}$  and false otherwise.

- Applications:

- File systems, web caching, databases, Google, checksum P2P network,...

# C++ unordered\_set

**cppreference.com**Create account

Page

Discussion

View

Edit

History

C++Containers librarystd::unordered\_set

## std::unordered\_set

Defined in header `<unordered_set>`

```
template<
    class Key,
    class Hash = std::hash<Key>,
    class KeyEqual = std::equal_to<Key>,
    class Allocator = std::allocator<Key>
> class unordered_set;

namespace pmr {
    template <class Key,
              class Hash = std::hash<Key>,
              class Pred = std::equal_to<Key>>
    using unordered_set = std::unordered_set<Key, Hash, Pred,
                                              std::pmr::polymorphic_allocator<Key>>;
}
```

(1) (since C++11)

(2) (since C++17)

Unordered set is an associative container that contains a set of unique objects of type Key. Search, insertion, and removal have average constant-time complexity.

# Dictionary problem

- Challenge:
  - The size of universe  $\mathbf{U}$  is too large to store in an array of size  $|\mathbf{U}|$ .
- Can we use the following data structures?
  - Bitvector
  - Linked list
  - Binary search tree
  - Hash table

# Hashing

- Hash function:

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

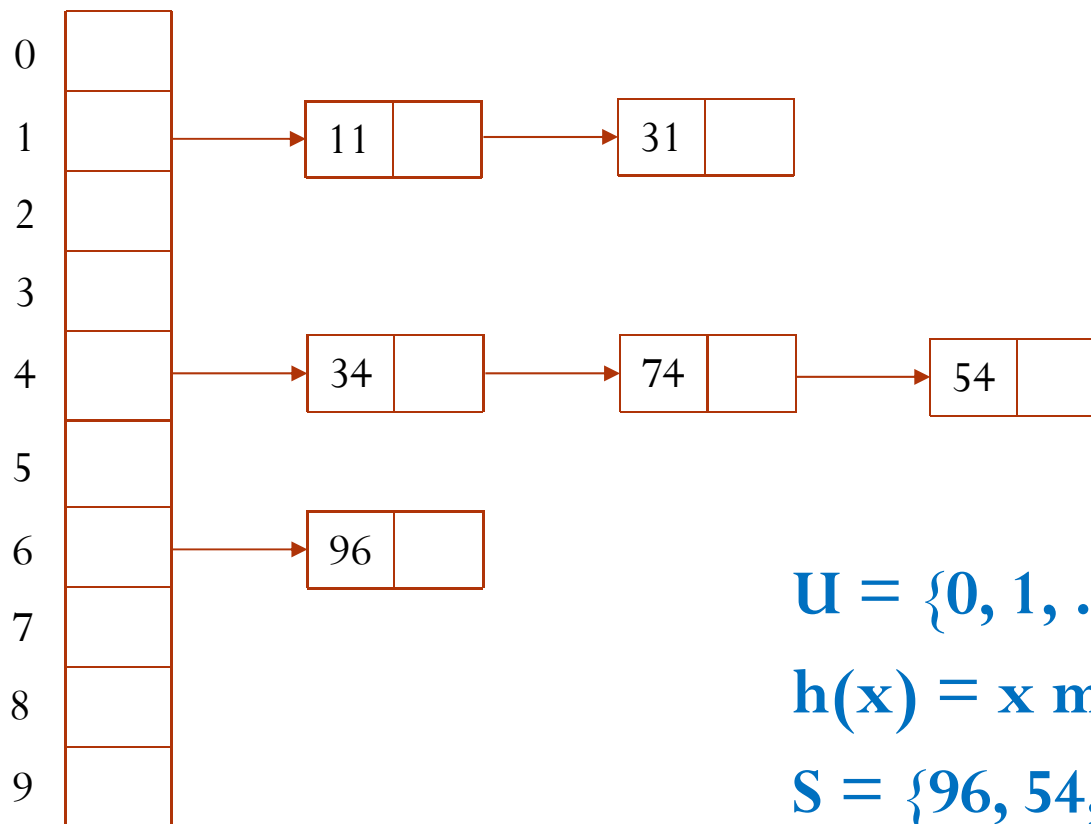
- Hashing:

- Create an array **A** of length **n**. For any operation on element  $x \in S$ , operate on the array position **A[h(x)]**

- Collision:

- Two different elements **x** and **y** have the same hash value:  $h(x) = h(y)$ .
- **Chaining**: **A[i]** is a linked list containing all elements **x** where  $h(x) = i$ .

# Example of chain hashing



$$U = \{0, 1, \dots, 99\}$$

$$h(x) = x \bmod 10$$

$$S = \{96, 54, 31, 11, 74, 34\}$$

# Dictionary operations

- **Search( $x$ ):**
  - Compute  $h(x)$ , then scan through the list  $A[h(x)]$ . Return true if finding  $x$ , and false otherwise.
- **Insert( $x$ ):**
  - Compute  $h(x)$ , then scan through the list  $A[h(x)]$  to find  $x$ . If  $x$  is not in the list, add  $x$  to the front of the list. Otherwise, do nothing.
- **Delete( $x$ ):**
  - Compute  $h(x)$ , then scan through the list  $A[h(x)]$  to find  $x$ . If  $x$  is in the list, remove  $x$  from the list. Otherwise, do nothing.
- **Time complexity:  $O(1 + \text{length of } A[h(x)])$**

# Deterministic hash function

- For any choice of  $h$ , we can find a set whose elements have the same hash values (map to one slot).
- We end up with a single linked list containing all elements. Hence, it takes  $O(n)$  time for the searching operation.
- Solutions:
  - Assume the input  $S$  is random.
  - **Ideal hash function:** Map  $m$  elements uniformly at random to  $n$  slots.



# Hashing performance

- Ideal hash function:
  - Any  $m$  elements are uniformly distributed to  $n$  slots.
  - Average length of chain =  $m/n$ .
  - Choose  $m \sim n$ , we expect  $O(1)$  cost for insert, delete, search operations.
- We cannot have an ideal hash function for all possible sets of elements (true randomness might not exist).
- Approach: Choose the hash function  $h$  at random from a hash family  $H$ .

# Universal hashing

- A universal family of hash functions  $\mathbf{H}$  is a **set** of hash functions  $\mathbf{h}$  mapping the universe  $\mathbf{U}$  to a set  $\{0, 1, \dots, n-1\}$  such that
    - For any different pair  $\mathbf{x} \neq \mathbf{y}$ :  $\Pr_{\mathbf{h}}[\mathbf{h}(\mathbf{x}) = \mathbf{h}(\mathbf{y})] \leq 1/n$ , i.e. collision probability is very tiny.
    - Can select a random  $\mathbf{h}$  efficiently
    - Can compute  $\mathbf{h}(\mathbf{x})$  efficiently (in constant time)
- $\mathbf{h}$  is chosen at random.

# Example

- Universal property: If  $\mathbf{x} \neq \mathbf{y}$ ,  $\Pr_h[\mathbf{h}(\mathbf{x}) = \mathbf{h}(\mathbf{y})] \leq 1/n$ .

Ex.  $U = \{a, b, c, d, e, f\}$ ,  $n = 2$ .

	a	b	c	d	e	f
$h_1(x)$	0	1	0	1	0	1
$h_2(x)$	0	0	0	1	1	1

	a	b	c	d	e	f
$h_1(x)$	0	1	0	1	0	1
$h_2(x)$	0	0	0	1	1	1
$h_3(x)$	0	0	1	0	1	1
$h_4(x)$	1	0	0	1	1	0

$$H = \{h_1, h_2\}$$

$$\Pr_h \bar{\cap}_H [h(a) = h(b)] = 1/2$$

$$\Pr_h \bar{\cap}_H [h(a) = h(c)] = 1$$

$$\Pr_h \bar{\cap}_H [h(a) = h(d)] = 0$$

...

$$H = \{h_1, h_2, h_3, h_4\}$$

$$\Pr_h \bar{\cap}_H [h(a) = h(b)] = 1/2$$

$$\Pr_h \bar{\cap}_H [h(a) = h(c)] = 1/2$$

$$\Pr_h \bar{\cap}_H [h(a) = h(d)] = 1/2$$

$$\Pr_h \bar{\cap}_H [h(a) = h(e)] = 1/2$$

$$\Pr_h \bar{\cap}_H [h(a) = h(f)] = 0$$

...

not universal

universal

# Universal hash family for integer

- We construct a universal hash function  $h$  that hashes an integer  $x \in U$  to a set of  $\{0, 1, \dots, n-1\}$  as follows:
  - Pick a large prime number  $p \geq n$
  - Pick a random pair of integers  $0 \leq a, b < p$
  - $h_{a,b}(x) = ((ax + b) \bmod p) \bmod n$
- Universal hash family:  $H = \{h_{a,b} : 0 \leq a, b < p\}$
- Proof of the universal property:
  - The 13.6 section of Algorithm Design textbook.
  - Carter & Wegman: “Universal Classes of Hash Functions”. STOC’77.

# Universal hash for a set of integers

- We construct a universal hash function  $\mathbf{h}$  that hashes a **set of integers**  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_d\}$  to a set  $\{0, 1, \dots, n-1\}$  as follows:
  - Choose independently at random  $d$  hash functions  $\mathbf{h}_1, \dots, \mathbf{h}_d$  from a universal hash family  $\mathbf{H}$ .
  - $\mathbf{h}(\mathbf{X}) = \mathbf{h}_1(\mathbf{x}_1) + \dots + \mathbf{h}_d(\mathbf{x}_d) \bmod n$ .
- **Proof of universal property:**
  - Carter & Wegman: “Universal Classes of Hash Functions”. STOC’77.
- **Practical hash function:** For a large prime  $p$  and  $0 \leq a_i < p$ ,
  - $\mathbf{h}_{a_0, \dots, a_d}(\mathbf{X}) = ((a_0 + a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \dots + a_d \mathbf{x}_d) \bmod p) \bmod n$

# Analysis of dictionary operations

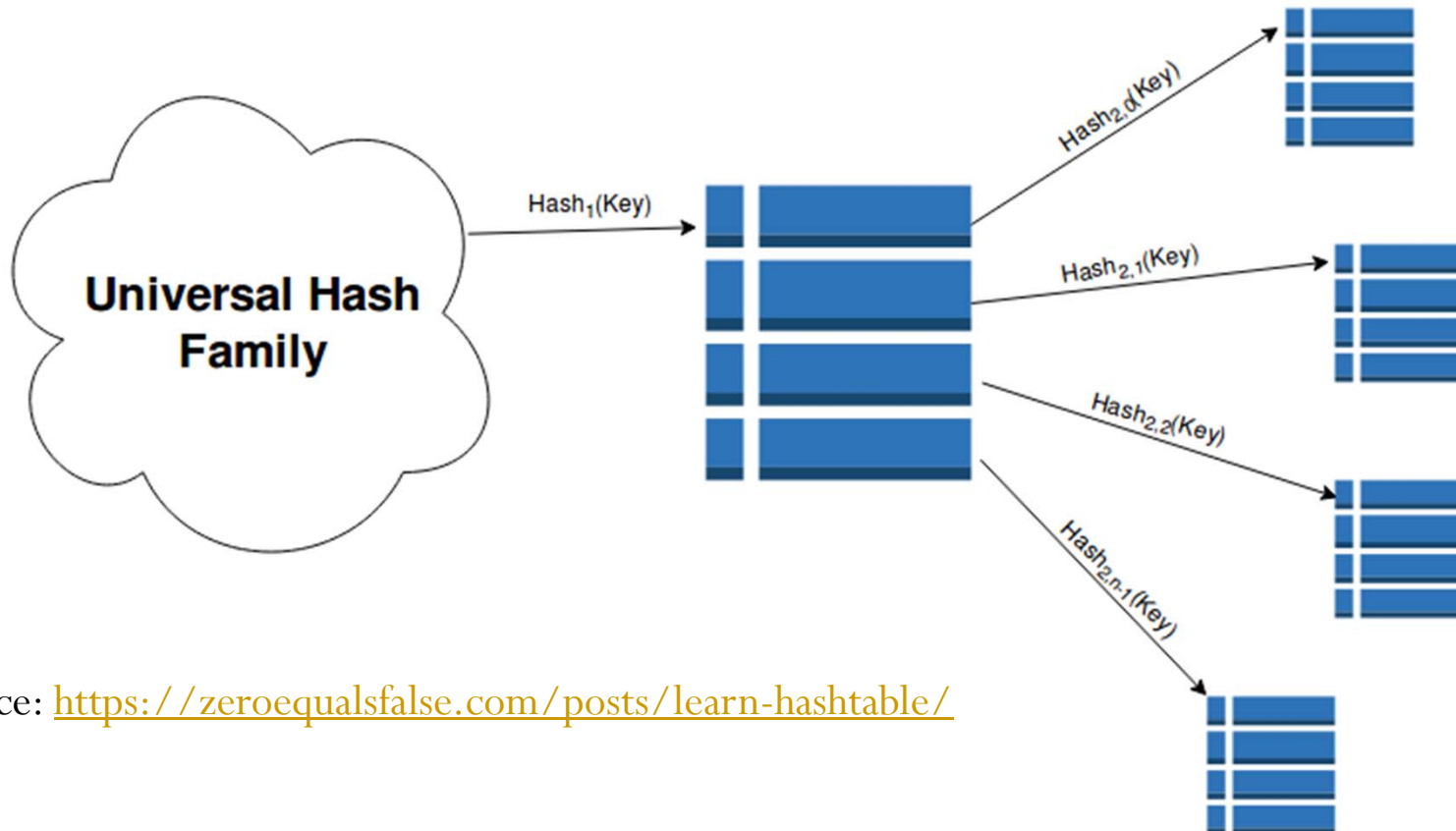
- Time complexity:  $O(1 + \text{length of } A[h(\mathbf{x})])$
- Assumption:  $|S| \leq n$  at all time, hence our dictionary uses  $O(n)$  space.
- The expectation of the length of the linked list  $A[h(\mathbf{x})]$ :
  - $$\begin{aligned} \mathbf{E} [\text{length of } A[h(\mathbf{x})]] &= 1 + \mathbf{E} [|\mathbf{y} \in S \setminus \{\mathbf{x}\} \mid h(\mathbf{y}) = h(\mathbf{x})|] \\ &\leq 1 + (|S| - 1) * \frac{1}{n} \leq 2. \end{aligned}$$
  - Expected time complexity of dictionary operations is  $O(1)$ .

# Conclusion

- Given a random choice of a hash function  $h$  from the universal family such that  $h: U \rightarrow \{0, 1, \dots, n-1\}$ .
- Our dictionary needs  $O(n)$  space and  $O(1)$  expected time per operation (search, insert, delete).
- Expectation is over the random choice of hash function.
- Independent of input set.

# Homework

- Write Python script to implement the universal hash family for a set of integers



Source: <https://zeroequalsfalse.com/posts/learn-hashtable/>