

# Algorithms

## (Hash Tables)

**Pramod Ganapathi**

Department of Computer Science  
State University of New York at Stony Brook

September 19, 2021



# Contents

- Hash Functions
- Hash Tables
- Collision-Avoiding Techniques
  - Separate Chaining
  - Open Addressing
- Applications

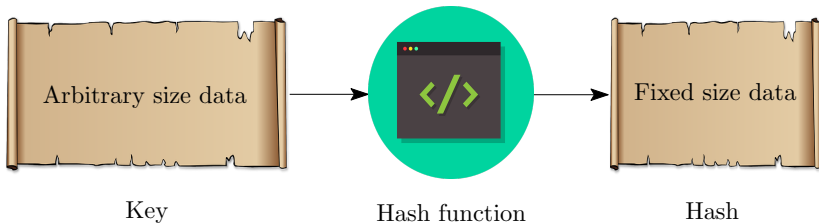
# Hash Functions

[HOME](#)

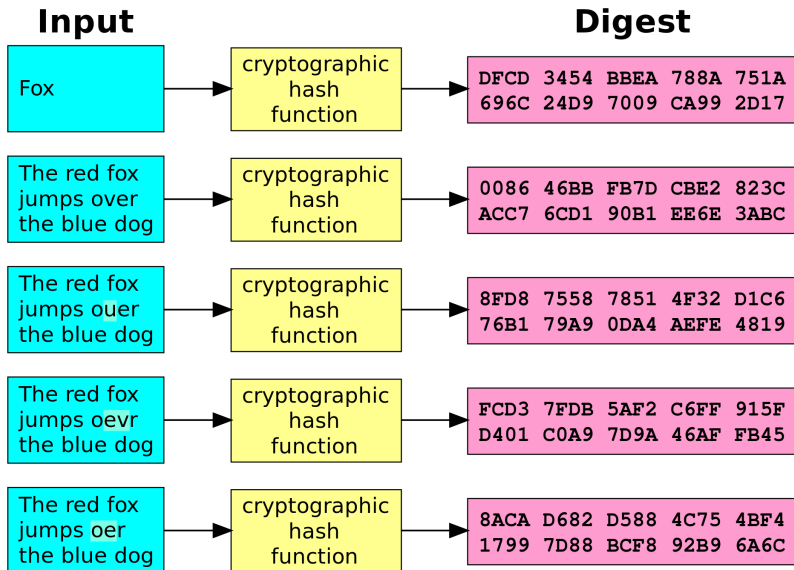
# Hash function

## Definition

- A hash function is a function that maps arbitrary size data to fixed size data.

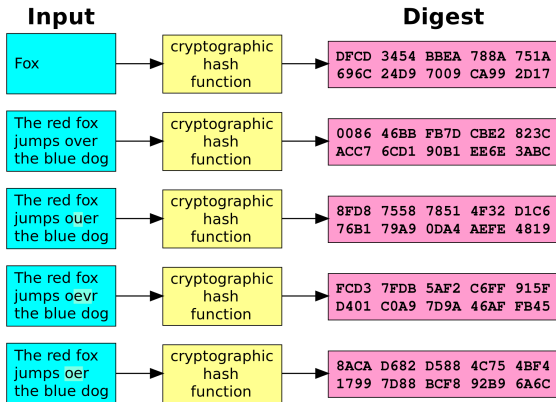


# Cryptographic hash function



Source: Wikipedia

# Properties of an ideal cryptographic hash function

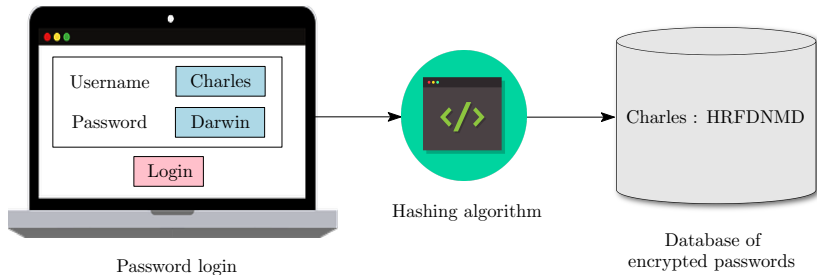


- Deterministic and fast
- Computing message from hash value is infeasible
- Computing two messages having same hash value is infeasible
- Tiny change in message must change the hash value drastically

# Applications of hashing

- Web page search using URLs
- Password verification
- Symbol tables in compilers
- Filename-filepath linking in operating systems
- Plagiarism detection using Rabin-Karp string matching algorithm
- English dictionary search
- Used as part of the following concepts:
  - finding distinct elements
  - counting frequencies of items
  - finding duplicates
  - message digests
  - commitment
  - Bloom filters

# Password authentication





# Hash Tables

[HOME](#)

# Hash table

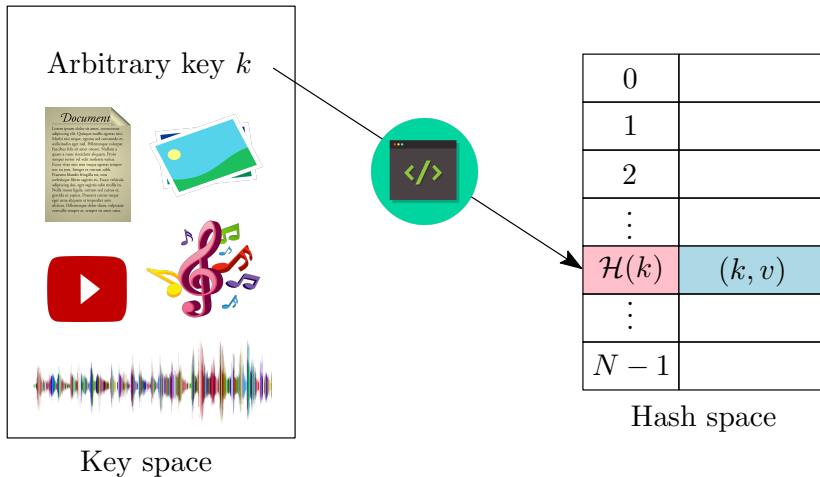
- A **hash table** is a data structure to implement **dictionary ADT**
- A hash table is an efficient implementation of a **set/multiset/map/multimap**.
- A hash table performs insert, delete, and search operations in **constant expected time**.

# Balanced search trees vs. Hash tables

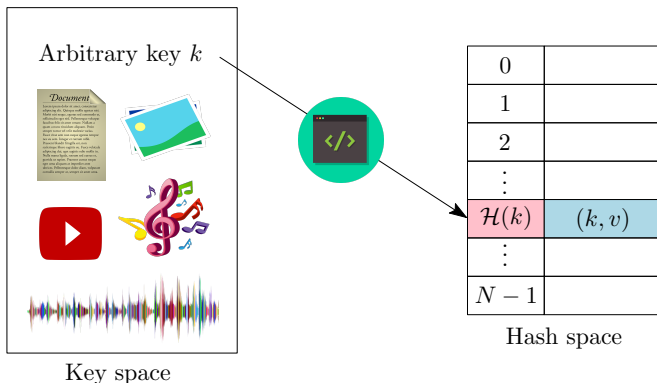
- Balanced search tree is ideal for sorted collection
- Hash table is ideal for unsorted collection

Operations		Balanced tree (worst)	Hash table (avg.)   (worst)	
Sorting-unrelated operations	Insert	$\mathcal{O}(\log n)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$
	Delete	$\mathcal{O}(\log n)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$
	Search	$\mathcal{O}(\log n)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$
Sorting-related operations	Sort	$\mathcal{O}(n)$	X	
	Minimum	$\mathcal{O}(\log n)$	X	
	Maximum	$\mathcal{O}(\log n)$	X	
	Predecessor	$\mathcal{O}(\log n)$	X	
	Successor	$\mathcal{O}(\log n)$	X	
	Range-Minimum	$\mathcal{O}(\log n)$	X	
	Range-Maximum	$\mathcal{O}(\log n)$	X	
	Range-Sum	$\mathcal{O}(n)$	X	

# Hash table



# Hash table



- A hash function is a mapping from arbitrary objects to the set of indices  $[0, N - 1]$ .
- The key-value pair  $(k, v)$  is stored at  $A[\mathcal{H}(k)]$  in the hash table.

# Hash table

## Questions

1. How can keys of arbitrary objects be mapped to array indices that are whole numbers?
2. How can an infinite number of keys be mapped to a finite number of indices?
3. What are the pros/cons of hash tables w.r.t. balanced trees?
4. What are the properties of an ideal hash function?
5. Is there one practical hash function that is best for all input?
6. What is the hash function used in Java?
7. Will there be collisions during hashing?  
If yes, how can we avoid collisions?
8. Is there a relation between table size and the number of items?
9. Why (key, value) pairs? Why not tuples?

# Encoding of information

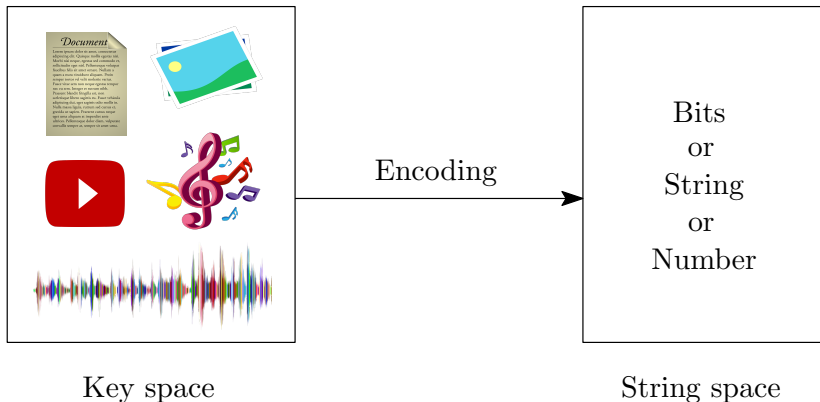
## Problem

- How can keys of arbitrary objects be mapped to array indices that are whole numbers?

# Encoding of information

## Problem

- How can keys of arbitrary objects be mapped to array indices that are whole numbers?





# Two stages of hash function

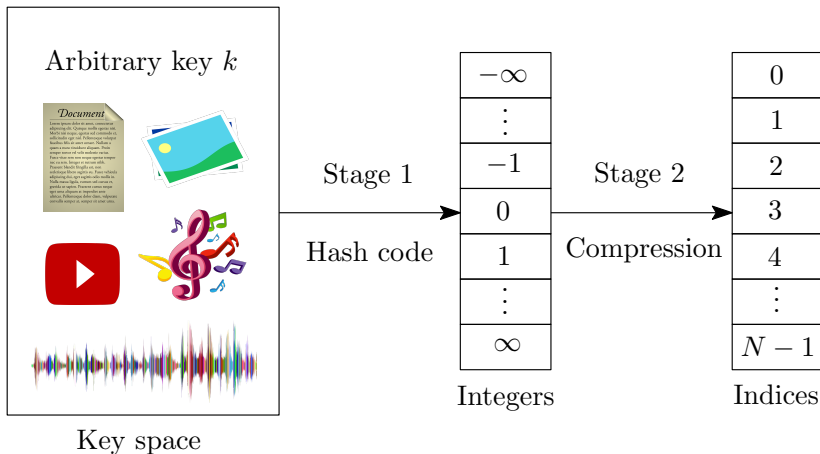
## Problem

- How can an infinite number of keys be mapped to a finite number of indices?

# Two stages of hash function

## Problem

- How can an infinite number of keys be mapped to a finite number of indices?



# Stage 1 of hash function: Hash code

- Consider bits as integer.

HashCode(byte | short | char) = 32-bit int ▷ upscaling

HashCode(float) = 32-bit int ▷ change representation

HashCode(double) = 32-bit int ▷ downscaling

HashCode( $x_0, x_1, \dots, x_{n-1}$ ) =  $x_0 + x_1 + \dots + x_{n-1}$  ▷ sum

HashCode( $x_0, x_1, \dots, x_{n-1}$ ) =  $x_0 \oplus x_1 \oplus \dots \oplus x_{n-1}$  ▷ xor

- Polynomial hash codes.

HashCode( $x_0, x_1, \dots, x_{n-1}$ ) =

$x_0 a^{n-1} + x_1 a^{n-2} + \dots + x_{n-2} a + x_{n-1}$  ▷ polynomial

- Cyclic-shift hash codes.

HashCode<sub>k</sub>( $x$ ) = Rotate( $x$ ,  $k$  bits) ▷ cyclic-shift

e.g.: Hashcode<sub>2</sub>(111000) = 100011

## Stage 2 of hash function: Compression function

A good compression function minimizes the number of collisions for a given set of distinct hash codes.

- **Division method.**

$$\text{Compression}(i) = i \% N$$

▷ remainder

$N \geq 1$  is the size of the bucket array.

Often,  $N$  being prime “spreads out” the distribution of primes.

Ex. 1: Insert codes  $\{200, 205, \dots, 600\}$  into  $N$ -sized array.

Which is better:  $N = 100$  or  $N = 101$ ?

Ex. 2: Insert multiple codes  $\{aN + b\}$  into  $N$ -sized array.

Which is better:  $N = \text{prime}$  or  $N = \text{non-prime}$ ?

- **Multiply-Add-and-Divide (MAD) method.**

$$\text{Compression}(i) = ((ai + b) \% p) \% N$$

▷ remainder

$N \geq 1$  is the size of the bucket array.

$p$  is a prime number larger than  $N$ .

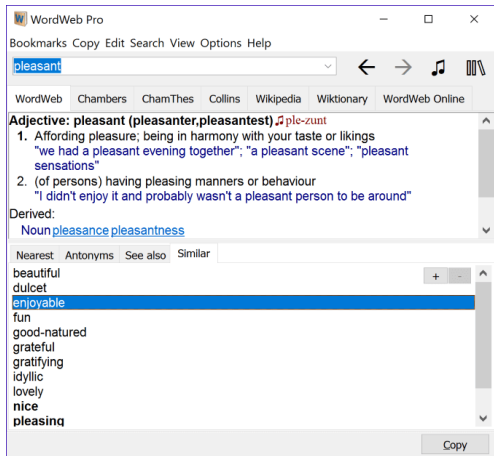
$a, b$  are random integers from the range  $[0, p - 1]$  with  $a > 0$ .

Usually eliminates repeated patterns in the set of hash codes.

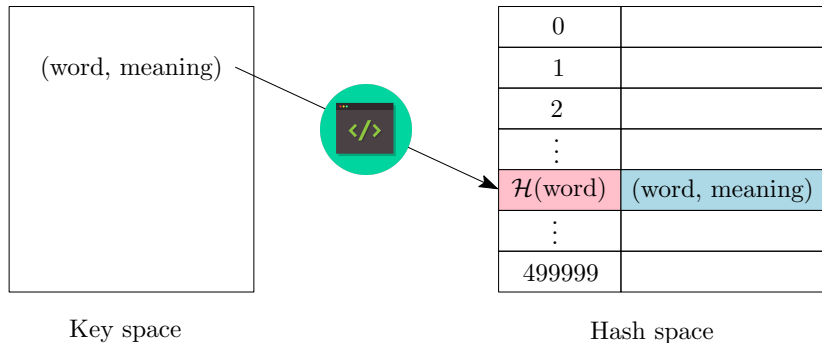
# English dictionary

## Problem

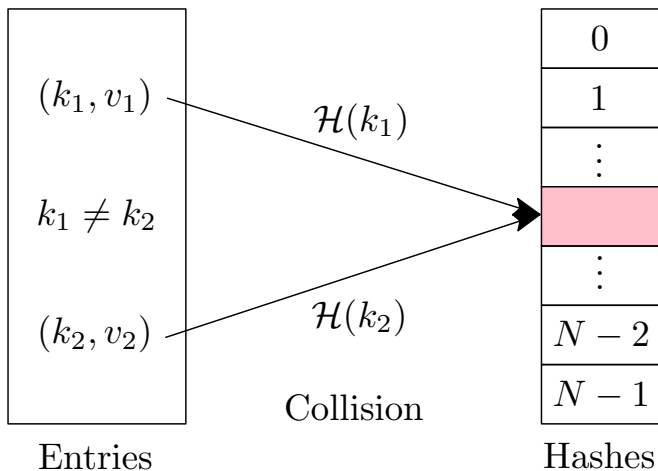
- How do you implement the English dictionary search such that searching for a word takes  $\mathcal{O}(1)$  time on an average?



# Hash table of English dictionary



# Collisions



# Collision-handling schemes

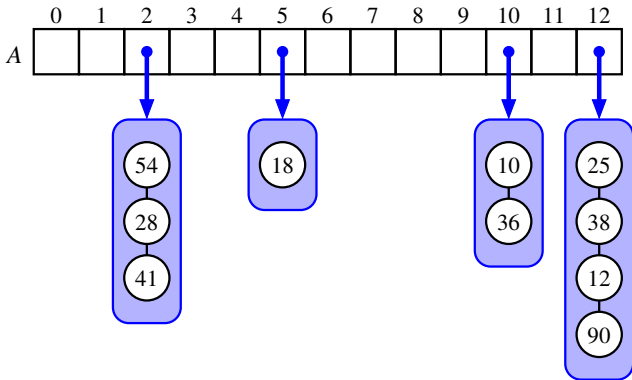
There are two major collision-handling schemes or collision-resolution strategies.

Collision-handling scheme	Features
Separate chaining	Extra space (for secondary data structures) Simpler implementation
Open addressing	No extra space More complicated implementation



# Separate chaining

- Have each bucket  $A[j]$  store its own secondary container.
- We use **secondary data structures** (e.g. array list, linked list, balanced search trees, etc) for each bucket.



# Separate chaining (via arraylist/linkedlist)

PUT( $(key, value)$ )

1.  $hash \leftarrow \text{HASH}(key)$
2.  $A[hash].\text{ADDLAST}((key, value))$   $\triangleright A[hash]$  is a linked list

GET( $key$ )

1.  $hash \leftarrow \text{HASH}(key)$
2. **return**  $A[hash].\text{GET}(key)$   $\triangleright$  returns value

REMOVE( $key$ )

1.  $hash \leftarrow \text{HASH}(key)$
2. **return**  $A[hash].\text{REMOVE}(key)$   $\triangleright$  returns removed value

# Open addressing

- All entries are stored in the bucket array itself.
- Strict requirement: **Load factor must be at most 1.**
- Useful in applications where there are space constraints, e.g.: smartphones and other small devices.
- Iteratively search the bucket  $A[(\text{HASH}(key) + f(i)) \% N]$  for  $i = 0, 1, 2, 3, \dots$  until finding an empty bucket.

Scheme	Function
Linear probing	$f(i) = i$
Quadratic probing	$f(i) = i^2$
Double hashing	$f(i) = i \cdot \text{HASH2}(key)$ e.g. $\text{HASH2}(key) = p - (key \% p)$ for prime $p < N$ . Here, $N$ should be a prime number.
Random generator	$f(i) = \text{RANDOM}(i, \text{HASH}(key))$

# Open addressing: Linear probing: Put

- Suppose  $\text{HASH}(\text{key}) = \text{key} \% 10$

Put			Array									
Key	→	Hash	0	1	2	3	4	5	6	7	8	9
18	→	8									18	
41	→	1		41							18	
22	→	2		41	22						18	
32	→	2		41	22						18	
(2 probes)				41	22	32					18	
98	→	8		41	22	32					18	
(2 probes)				41	22	32					18	98
58	→	8		41	22	32					18	98
(3 probes)				41	22	32					18	98
78	→	8	58	41	22	32					18	98

How many probes are required to insert 78?

# Open addressing: Linear probing: Remove

- Suppose  $\text{HASH}(\text{key}) = \text{key} \% 10$

Remove Key	Array									
	0	1	2	3	4	5	6	7	8	9
—	58	41	22	32	78	19			18	98
58	58	41	22	32	78	19			18	98
		41	22	32	78	19			18	98
19		41	22	32	78	19			18	98

Hence, we cannot simply remove a found entry.

Remove Key	Array									
	0	1	2	3	4	5	6	7	8	9
—	58	41	22	32	78	19			18	98
58	58	41	22	32	78	19			18	98
	58	41	22	32	78	19			18	98
19	58	41	22	32	78	19			18	98
	58	41	22	32	78	19			18	98

Replace the deleted entry with the defunct object.

# Open addressing: Linear probing

PUT(*(key, value)*)

1.  $hash \leftarrow \text{HASH}(key); i \leftarrow 0$
2. **while**  $(hash + i) \% N \neq \text{null}$  **and**  $i < N$  **do**  $i \leftarrow i + 1$
3. **if**  $i = N$  **then throw** Bucket array is full
4. **else**  $A[(hash + i) \% N] \leftarrow (key, value)$

GET(*key*)

1.  $hash \leftarrow \text{HASH}(key); i \leftarrow 0$
2. **while**  $(hash + i) \% N \neq \text{null}$  **and**  $i < N$  **do**
3.    $index \leftarrow (hash + i) \% N$
4.   **if**  $A[index].key = key$  **then return**  $A[index].value$
5.    $i \leftarrow i + 1$
6. **return null**

REMOVE(*key*)

1.  $index \leftarrow \text{FINDSLOTFORREMOVAL}(key)$
2. **if**  $index < 0$  **then return null**
3.  $value \leftarrow A[index].value; A[index] \leftarrow \text{defunct}; n \leftarrow n - 1$
4. **return value**

# Complexity

- Suppose  $N$  = bucket array size and  $n$  = number of entries.
- Ratio  $\lambda = n/N$  is called the **load factor** of the hash table.
- If  $\lambda > 1$ , **rehash**. Make sure  $\lambda < 1$ .
- Assuming good hash function, expected size of bucket is  $\mathcal{O}(\lceil \lambda \rceil)$ .
- Separate chaining: Maintain  $\lambda < 0.75$   
Open addressing: Maintain  $\lambda < 0.5$
- Assuming good hash function and  $\lambda \in \mathcal{O}(1)$ ,  
complexity of **put**, **get**, and **remove** is  $\mathcal{O}(1)$  expected time.

# Applications

HOME



# Symbol tables in compilers

```
1. double foo(int count)
2. {
3.     double sum = 0.0;
4.     for (int i = 1; i <= count; i++)
5.         sum += i;
6.     return sum;
7. }
```

Symbol	Type	Scope
foo	function, double	global
count	int	function parameter
sum	double	block local
i	int	for-loop statement