

# Lecture 4: Hash Tables

김수경

이화여자대학교 인공지능융합전공 소속

# 저작권 안내

---

(주)업스테이지가 제공하는 모든 교육 콘텐츠의 지식재산권은  
운영 주체인 (주)업스테이지 또는 해당 저작물의 적법한 관리자에게 귀속되어 있습니다.

콘텐츠 일부 또는 전부를 복사, 복제, 판매, 재판매 공개, 공유 등을 할 수 없습니다.

유출될 경우 지식재산권 침해에 대한 책임을 부담할 수 있습니다.

유출에 해당하여 금지되는 행위의 예시는 다음과 같습니다.

- 콘텐츠를 재가공하여 온/오프라인으로 공개하는 행위
- 콘텐츠의 일부 또는 전부를 이용하여 인쇄물을 만드는 행위
- 콘텐츠의 전부 또는 일부를 녹취 또는 녹화하거나 녹취록을 작성하는 행위
- 콘텐츠의 전부 또는 일부를 스크린 캡처하거나 카메라로 촬영하는 행위
- 지인을 포함한 제3자에게 콘텐츠의 일부 또는 전부를 공유하는 행위
- 다른 정보와 결합하여 Upstage Education의 콘텐츠를 알 수 있는 저작물을 작성, 공개하는 행위
- 제공된 데이터의 일부 혹은 전부를 Upstage Education 프로젝트/실습 수행 이외의 목적으로 사용하는 행위

## Data Structures So Far

	Insertion	Retrieval	Deletion	Restriction
Sorted Array	$O(M)$	$O(\log M)$	$O(M)$	
Sorted Linked list	$O(M)$	$O(M)$	$O(M)$	
Stack	$O(1)$	$O(1)$	$O(1)$	LIFO
Queue	$O(1)$	$O(1)$	$O(1)$	FIFO
Hash Table	$O(1)$	$O(1)$	$O(1)$	

Needed for applications that  
need radically fast operations:

- 911 emergency calls and locating caller's address
- Airline information system
- Web search

01

# Data Indexed Arrays

# Data Indexed Arrays

- A regular list: [2, 5, 9, 10]
  - In this setting, we take  $O(N)$  to search a particular value.

Index	[0]	[1]	[2]	[3]
Value	2	5	9	10

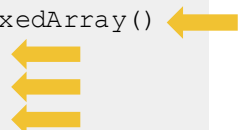
- What if we represent the data in a **data-indexed array**?

Index	[0]	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]	[9]	[10]
Value	F	F	T	F	F	T	F	F	F	T	T

# Data Indexed Arrays

- A data-indexed array has **all possible data as its indices!**
  - Initially, all values of the array are **False** (*i.e.*, `di_array[x] = False` for all `x` in `di_array`), meaning the array is empty.
  - When we **insert** some value, the corresponding index becomes **True**.
  - Here, we assume that no duplicate keys are allowed.

```
di_array = DataIndexedArray()
di_array.insert(3)
di_array.insert(6)
di_array.delete(3)
```



Index	[0]	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]	[9]	[10]
Value	F	F	F	T	F	F	T	F	F	F	F

# Data Indexed Arrays

- Why is this good?
  - We can insert, delete, and retrieve values in  **$O(1)$** !
- Any problem then?
  - We can't deal with values larger than 10 (the array size).
  - In other words, we need to use **very large memory space** if we'd like to deal with large number, even though the data size itself is small.
  - It is **not trivial how to deal with non-integer values**, like string or floating point numbers.

```
di_array.insert(21057381) 🤯
```

```
di_array.insert("ewha") 🤯
```

# Strings in Data Indexed Arrays

- Recall that each character can be represented as an ASCII code value (0~255).

○ “ewha”:

■ e: 101

■ w: 87

■ h: 104

■ a: 97

○ “ewha”<sub>256</sub>:  $(101 \times 256^3) + (87 \times 256^2) + (104 \times 256) + 97$

- In this way, a string can be represented as a **unique integer**.

DEC	ASCII	DEC	ASCII	DEC	ASCII	DEC	ASCII	DEC	ASCII	DEC	ASCII	DEC	ASCII	DEC	ASCII
1	!	32	space	64	@	96	`	128	Ç	160	à	192	À	224	Ô
2	"	33	!	65	A	97	a	129	à	161	á	193	Á	225	Õ
3	"	34	"	66	B	98	b	130	â	162	â	194	Â	226	Ö
4	+	35	#	67	C	99	c	131	ã	163	ã	195	Ã	227	Ø
5	*	36	\$	68	D	100	d	132	ä	164	ä	196	Ä	228	Ù
6	*	37	%	69	E	101	e	133	å	165	Å	197	Å	229	Ú
7	*	38	&	70	F	102	f	134	ä	166	ä	198	Ä	230	Û
8	u	39	'	71	G	103	g	135	ç	167	ç	199	Ç	231	Ü
9	o	40	(	72	H	104	h	136	è	168	è	200	È	232	Ý
10	m	41	)	73	I	105	i	137	é	169	é	201	É	233	Û
11	o	42	*	74	J	106	j	138	ê	170	ê	202	Ê	234	Ü
12	o	43	+	75	K	107	k	139	ë	171	ë	203	Ë	235	Ü
13	Z	44	,	76	L	108	l	140	ì	172	ì	204	Ì	236	ý
14	o	45	-	77	M	109	m	141	í	173	í	205	Í	237	ÿ
15	o	46	.	78	N	110	n	142	ä	174	ä	206	Ä	238	-
16	o	47	/	79	O	111	o	143	Å	175	Å	207	å	239	-
17	o	48	0	80	P	112	p	144	É	176	É	208	é	240	-
18	o	49	1	81	Q	113	q	145	æ	177	æ	209	æ	241	-
19	o	50	2	82	R	114	r	146	Æ	178	Æ	210	Æ	242	-
20	o	51	3	83	S	115	s	147	ó	179	ó	211	Ó	243	-
21	o	52	4	84	T	116	t	148	ô	180	ô	212	Ô	244	-
22	o	53	5	85	U	117	u	149	ö	181	ö	213	Ö	245	-
23	o	54	6	86	V	118	v	150	û	182	À	214	À	246	-
24	o	55	7	87	W	119	w	151	ü	183	Ä	215	Ä	247	-
25	o	56	8	88	X	120	x	152	ÿ	184	Ø	216	Ø	248	-
26	o	57	9	89	Y	121	y	153	Ö	185	Ø	217	Ø	249	-
27	o	58	:	90	Z	122	z	154	Ü	186	Ø	218	Ø	250	-
28	o	59	;	91	[	123	{	155	ø	187	Ø	219	Ø	251	-
29	o	60	<	92	\	124		156	£	188	Ø	220	Ø	252	-
30	o	61	=	93	]	125	}	157	Ø	189	Ø	221	Ø	253	-
31	o	62	>	94	^	126	~	158	×	190	Ø	222	Ø	254	-
		63	?	95	_	127	o	159	f	191	Ø	223	Ø	255	space

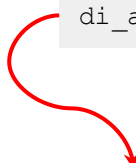
[1]



# Large Numbers in Data Indexed Arrays

- But still, we can't avoid the problem of treating large numbers.
  - We may need some way to put **arbitrary integers** in a **reasonable size of memory**.

```
di_array.insert(21057381)
```



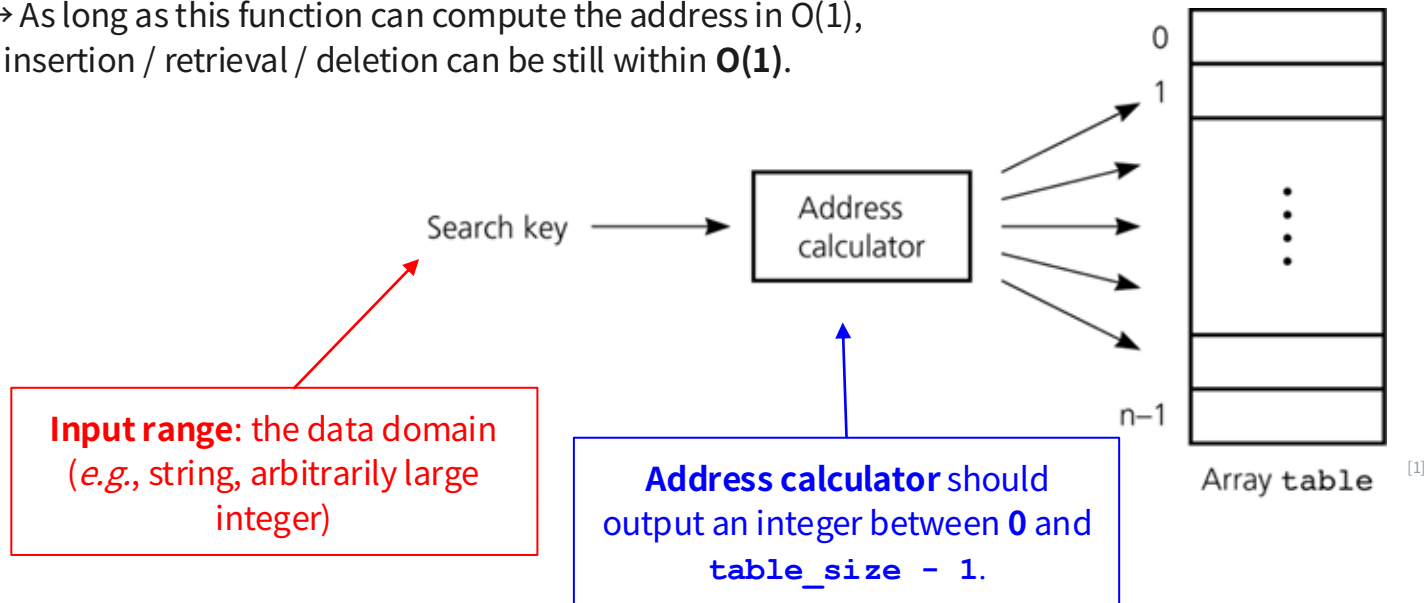
Index	[0]	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]	...	[999]
Value	F	F	F	F	F	F	F	F	F	...	F

# Hash Tables

# Hash Table

Given a key, a function (called **hash function**) determines where to locate it.

→ As long as this function can compute the address in  $O(1)$ , insertion / retrieval / deletion can be still within  **$O(1)$** .



[1] <https://www.slideshare.net/slideshow/hashing-in-data-structure-is-presented-in-these-slides/274382978>

# Hash Functions

- A good hash function:
  - Easy and **fast** to compute (in  $O(1)$ )
  - **Scatter** the data **evenly** on the hash table.
- Selection digits
  - $h(001364825) = 35$
- Folding
  - $h(001364825) = 001 + 364 + 825 = 1190$
- Modulo arithmetic
  - $h(x) = x \bmod \text{table\_size}$
  - $h(1004) = 4$  (if  $\text{table\_size} = 100$ )
  - This is widely used, as it maps the integers **uniformly** over the internal array regardless of its size!

Hmm, then what should we do if two different keys coincidentally **map to the same address**?

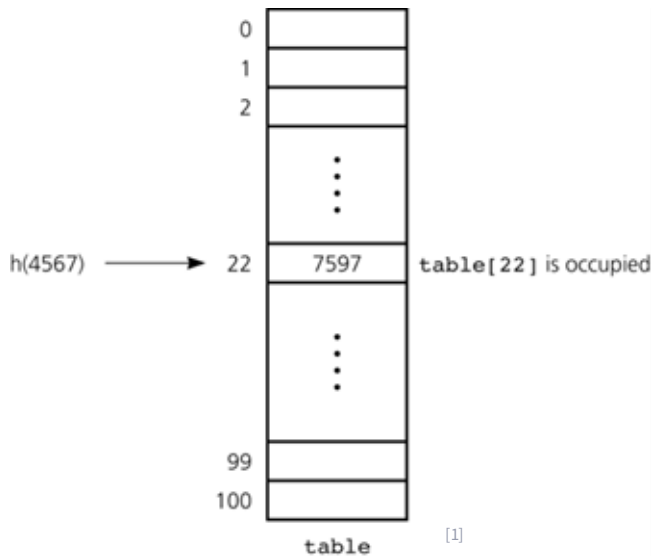


## Modulo operator (%)

$x \% y$  is the remainder when we divide  $x$  by  $y$ .

# Collision

- Yes, unfortunately, it can happen with all hash functions we listed 😬.
  - **Collision:** the phenomenon that two keys are mapped into the same location in the hash table.
- Selection digits
  - $h(001364825) = h(741325385) = 35$
- Folding
  - $h(001364825) = h(825364001) = 1190$
- Modulo arithmetic
  - $h(2205) = h(2405) = 5$  if `table_size = 100`



[1] <https://www.cs.colostate.edu/~cs165/Spring20/slides/16-Hashing.pdf>

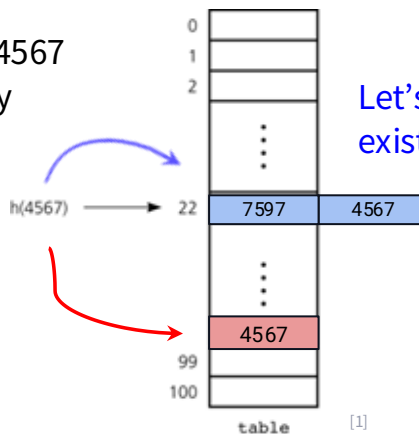
# Collision Resolution

- Collision is unavoidable; so, we need a method to resolve it efficiently.
- There are two groups of approaches:
  - **Open addressing**: allowing the items to be located in another place if the collision happened.
  - **Separate chaining**: allowing more than one item to be located at a place.

Example: According to our hash function, 4567 needs to be located at [22], but it is already occupied 😬.

## Open Addressing

Let's find another place to put it.



## Separate chaining

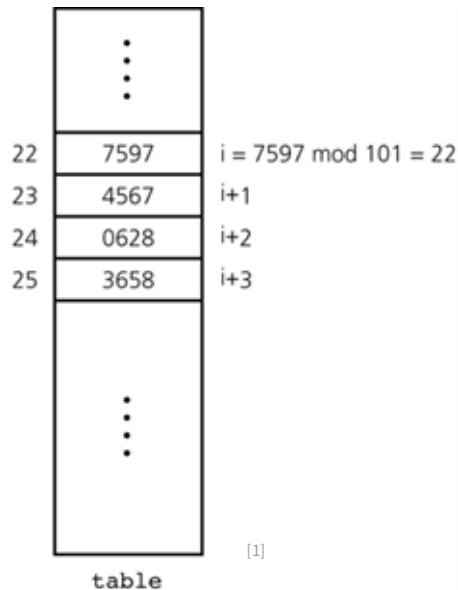
Let's put it there together with the existing one(s).

# Open Addressing

- **Linear probing:** if the designated spot is occupied, try the **right next index**, until it finds an available one.

$$h_i(x) = (h_0(x) + i) \bmod \text{table\_size}$$

- Causes **primary clustering**:
  - All keys mapped to [22], [23], [24], [25], ... will suffer from collision.
  - If this cluster gets bigger, insertion / retrieval / deletion of keys corresponding to this region will take  $O(N)$  instead of  $O(1)$ .
  - Gets significantly worse especially if consecutive numbers are sequentially added. (e.g.,  $22 \rightarrow 23 \rightarrow 24, \dots$ )



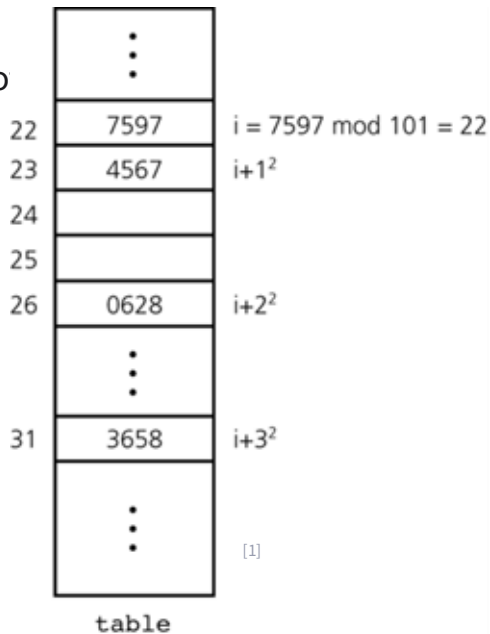
[1] <https://www.cs.colostate.edu/~cs165/Spring20/slides/16-Hashing.pdf>

# Open Addressing

- **Quadratic probing:** if the designated spot is occupied, try the next set of indices, **quadratically increasing** with the number of trials.

$$h_i(x) = (h_0(x) + i^2) \bmod \text{table\_size}$$

- Solves **primary clustering**:
  - Items collided at [23] suffers less from significant collisions originally happened at [22].
- Still suffers from **secondary clustering**:
  - For items collided at the same index for the first time, they must always follow the same trace.



[1] <https://www.cs.colostate.edu/~cs165/Spring20/slides/16-Hashing.pdf>



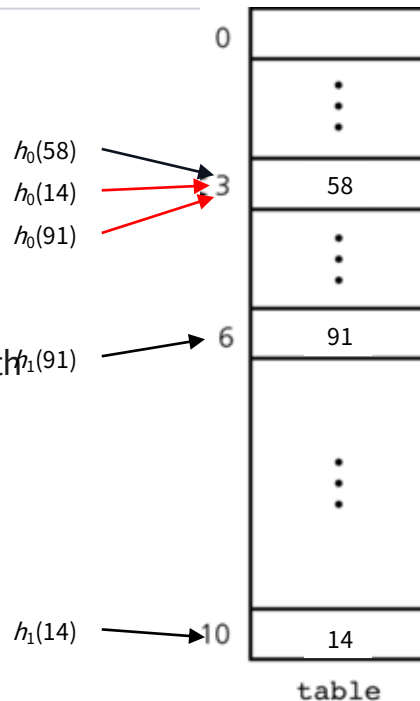
# Open Addressing

- **Double Hashing:** using two different hash functions to reduce collision probability.

$$h_i(x) = (\alpha(x) + i * \beta(x)) \bmod \text{table\_size}$$

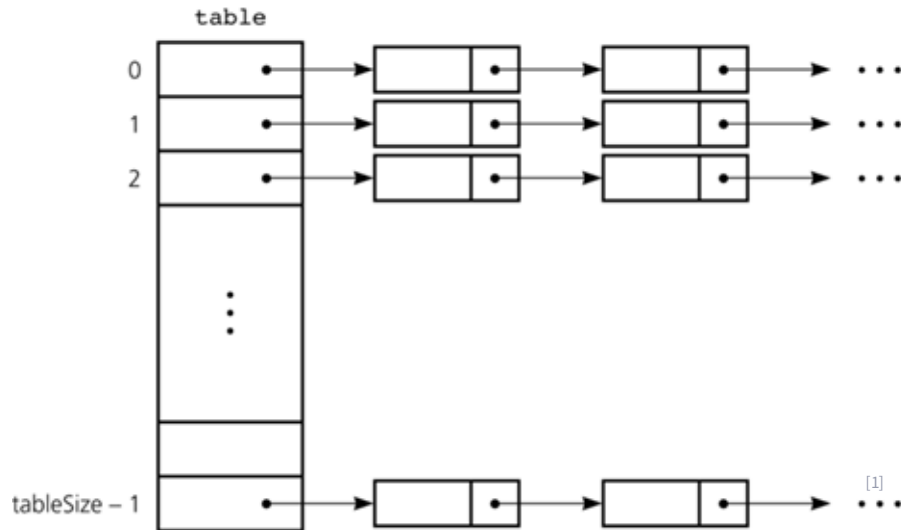
where  $\alpha(x), \beta(x)$ : some hash functions

- Can avoid primary and secondary clustering, since we probe with  $h_1^{(91)}$  different steps depending on the input.



# Separate Chaining

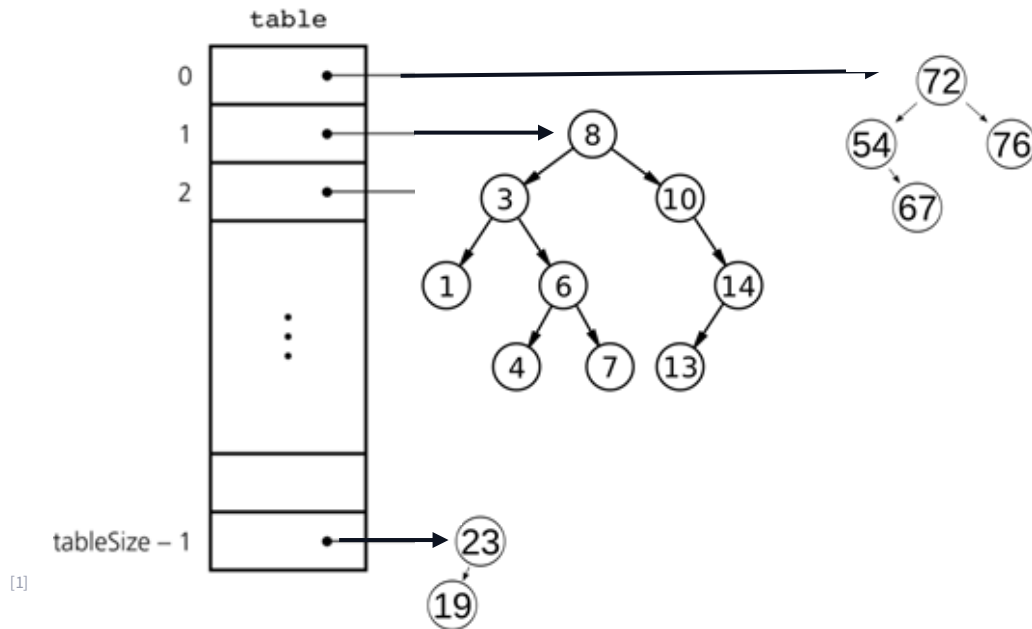
- Each location of the hash table is allowed to have more than one item.
  - A simplest example: a reference to a linked list.



[1] <https://medium.com/@ramyjbh/data-structures-for-dummies-hash-tables-579ddd1a4389>

# Separate Chaining

- Actually, we can put any other data structure that support key-based search.

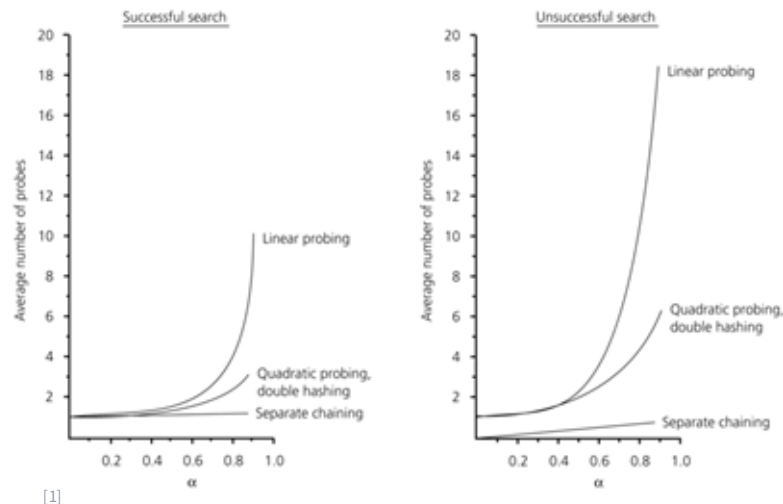


# Performance of a Hash Table

---

- In general, having more items in the table reduces performance of a hash table.
  - This is unique to hash tables, unlike most other data structures!
- Let's define **load factor ( $\alpha$ )** formally:  $N/M$ 
  - $N$ : the actual number of elements inserted
  - $M$ : the number of indices (table\_size)
- With **open addressing**,  $\alpha$  is always between 0 and 1. Once it reaches to 1, we can't add more data because there is no more available space.
  - If  $\alpha$  gets **close to 1**, there are probably lots of collisions, meaning **longer processing time**.
- With **separate chaining**, load factor can be  $> 1$ .
  - In this case, load factor means the average chain length per index, which is equivalent to the **expected number of linear traversal** to locate an item.
  - Obviously, higher load factor means slower data processing.

# Efficiency of Hashing



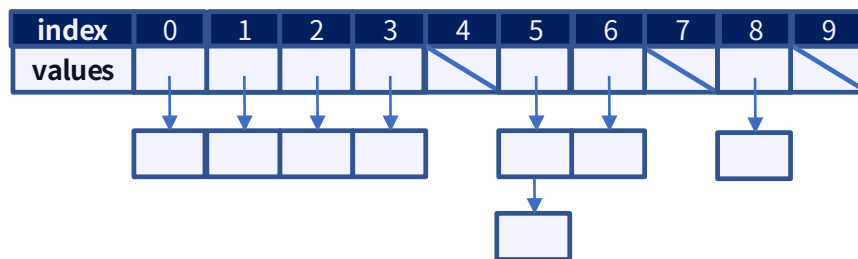
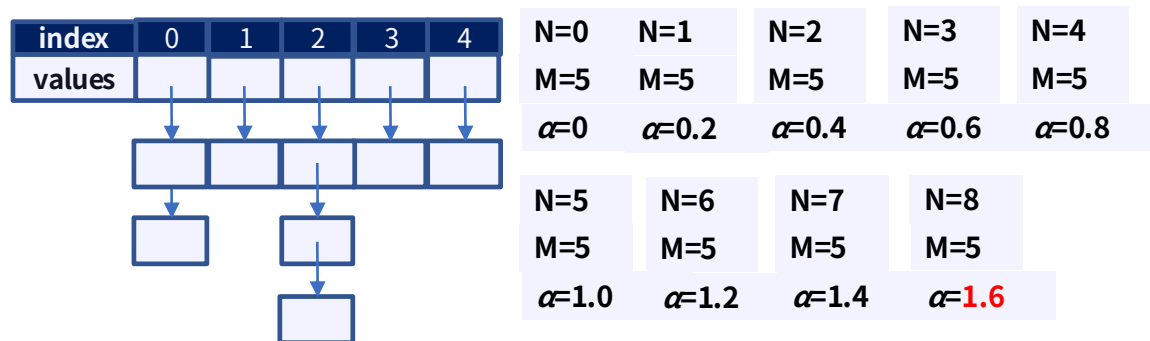
When the load factor is low, hash table achieves  $\approx O(1)$ , regardless of the collision resolution method.

## How to Decide the Table Size ( $M$ )

---

- If we know the data size ( $N$ ) to be inserted, we can set  $M = N/\alpha$ , where  $\alpha$  is the desired load factor.
  - With small  $\alpha$ , we focus more on speed.
  - With large  $\alpha$ , we care more about the memory space.
  - It is known that  $\alpha = 0.6 \sim 0.75$  is a good balance between them.
- What if we don't know the data size?
  - We have no choice but to dynamically adapt.
  - Starting with some default size (say, 101), then (roughly) double it when the load factor becomes larger than some threshold (e.g., 0.6~0.75).

# Resizing Example

**Resizing!****Redistributing all items!**

N=8  
M=10  
 $\alpha=0.8$

# Time Complexity of Resizing

---

- Besides resizing, hash table achieves  $O(1)$  insertion, deletion, and retrieval.
- Resizing is obviously not free 😓
  - Resizing a hash table with  $N$  items requires  $O(N)$  time to redistribute all  $N$  items.
- A good news 😊: we don't always resize since one resizing operation doubles the number of indices.
  - The number of redistributing items while inserting  $N$  items:
    - $1 + 2 + 4 + 8 + \dots + N = 2N - 1$
  - Overall, redistributing cost becomes  $O((2N-1)/N) = O(1)$  on average!



03

# Implementations: Dictionary vs Set in Python

# Python Dictionary

- Python internally provides the hash table under the name of Dictionary.

- **Key-value pairs** are stored in a hash table.
- Key are unique and hashable.
- Insertion, retrieval, deletion supported in  **$O(1)$** .

```
>>> tel = {'jerry': 1086, 'jose': 8249}
>>> tel['soo'] = 4127
>>> tel
{'jerry': 1086, 'jose': 8249, 'soo': 5564}
>>> tel['jose']
8249
>>> del tel['jerry']
>>> tel['shawn'] = 8080
>>> tel
{'jose': 8249, 'soo': 5564, 'shawn': 8080}
>>> list(tel)
['jose', 'soo', 'shawn']
>>> sorted(tel)
['shawn', 'soo', 'jose']
>>> 'jinri' in tel
False
>>> 'jeongwoo' not in tel
True
```

Insertion

Retrieval

Deletion

# Python Set

- If you don't have an associated value, you may use Python Set, instead of Dictionary.
  - **Only elements (keys)** are stored in the hash table.
  - Elements are unique and hashable.
  - Insertion, retrieval, deletion supported in **O(1)**.

```
>>> a = set('abracadabra')
>>> b = set('alacazam')
>>> a
{'a', 'r', 'b', 'c', 'd'}    Unique letters in a
>>> a - b
{'r', 'd', 'b'}    Set difference
>>> a | b
{'a', 'c', 'r', 'd', 'b', 'm', 'z', 'l'}    Set OR
>>> a & b
{'a', 'c'}    Set AND
>>> a ^ b
{'r', 'd', 'b', 'm', 'z', 'l'}    Set XOR
```

# Applications of Hash Tables

# Problem 1 - Palindrome

## 409. Longest Palindrome

Easy

Topics

Companies

Given a string `s` which consists of lowercase or uppercase letters, return the length of the **longest palindrome** that can be built with those letters.

Letters are **case sensitive**, for example, `"Aa"` is not considered a palindrome.

### Example 1:

Input: `s = "abcccccdd"`

Output: 7

Explanation: One longest palindrome that can be built is `"dccaccd"`, whose length is 7.

### Example 2:

Input: `s = "a"`

Output: 1

Explanation: The longest palindrome that can be built is `"a"`, whose length is 1.

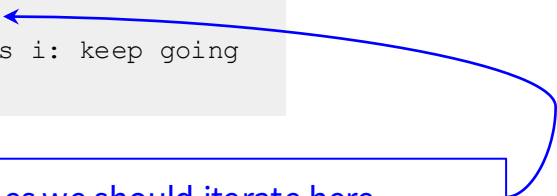
## Problem 2 - Smallest Missing Integer

- Given a list of integers, return the smallest positive integer which is not in that list.

Examples:

- [7, 2, 3, 5, 4, 1] → 6
- [-1, 5, 2, 3, 9] → 1
- [17, 25, 4308, 1, 99] → 2

```
# TODO(students): implement this!  
def smallest_missing_pos_int(self, list):  
    for each item in the list:  
        insert into a hash table (hastset)  
  
    for i = 1, 2, 3, ...:  
        if the hash table contains i: keep going  
        else: return i
```



**Note:** think about how many times we should iterate here.

# Problem 3 – Two Sum

## 1. Two Sum

Solved Easy 🔖 Topics 🏢 Companies 💡 Hint

Given an array of integers `nums` and an integer `target`, return *indices of the two numbers such that they add up to* `target`.

You may assume that each input would have **exactly one solution**, and you may not use the same element twice.

You can return the answer in any order.

### Example 1:

Input: `nums = [2,7,11,15]`, `target = 9`

Output: `[0,1]`

Explanation: Because `nums[0] + nums[1] == 9`, we return `[0, 1]`.

### Example 2:

Input: `nums = [3,2,4]`, `target = 6`

Output: `[1,2]`

### Example 3:

Input: `nums = [3,3]`, `target = 6`

Output: `[0,1]`



# Building intelligence for the future of work