

Algorithms & Data Structures I:

Hashing

Today's Topics

- Dictionaries
- Motivation
- Hashing
- Chaining
- Simple Uniform hashing
- «Good» hash functions

Dictionary (Map in C++)

Dictionary is an:

- Abstract Data Type (ADT) maintaining items, where each item is a pair<key, value>

Examples:

1. *Phonebook*. Keys are names, and their corresponding items are phone numbers
2. *Real dictionary*. Keys are english words, and their corresponding items are dictionary-entries

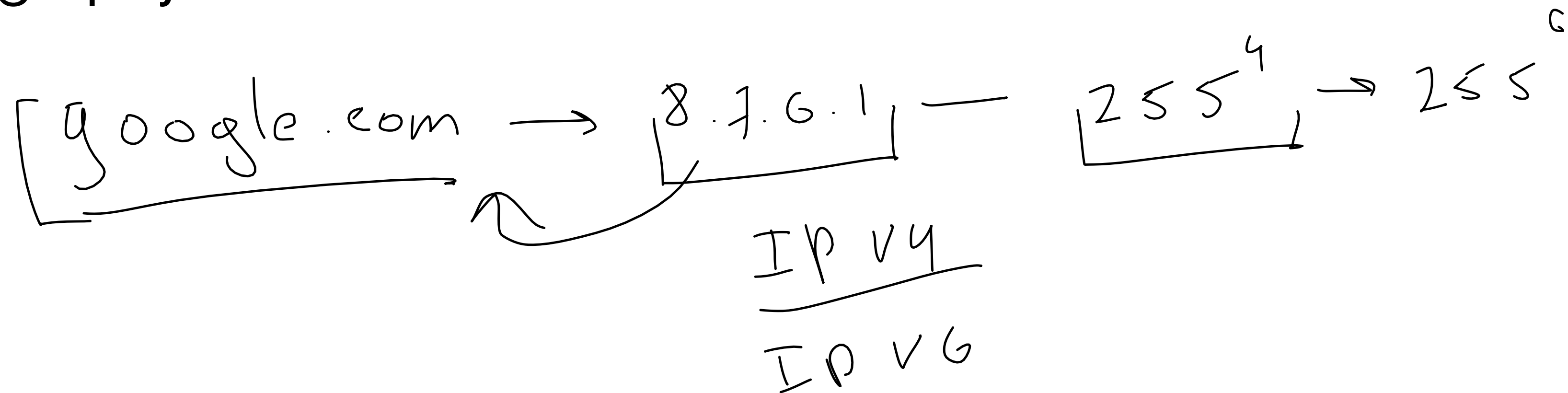
key : value
<James : +7-825-...>
<Johny : +7-7-7-...>

Real Dictionary:
key : value
ball :

Motivation

Dictionaries:

- built into most modern programming languages (Python, C++, Ruby, Go, JavaScript, Java, ...)
- very powerful concept
- use in web development fundamentals such as DNS system
- use in cryptography



Operations to support

Insert(item): Add item to the data structure

Insert("a", 1)

Insert("a", 2)

Delete(item): Delete item from the data structure

↳ "a", 2

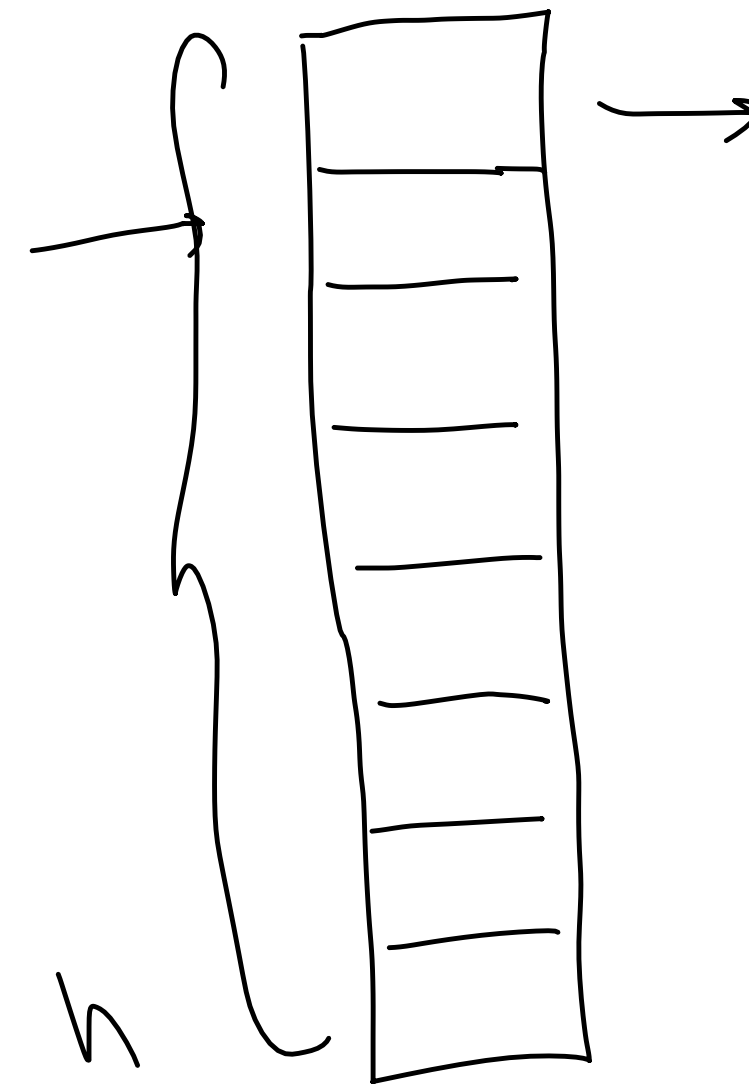
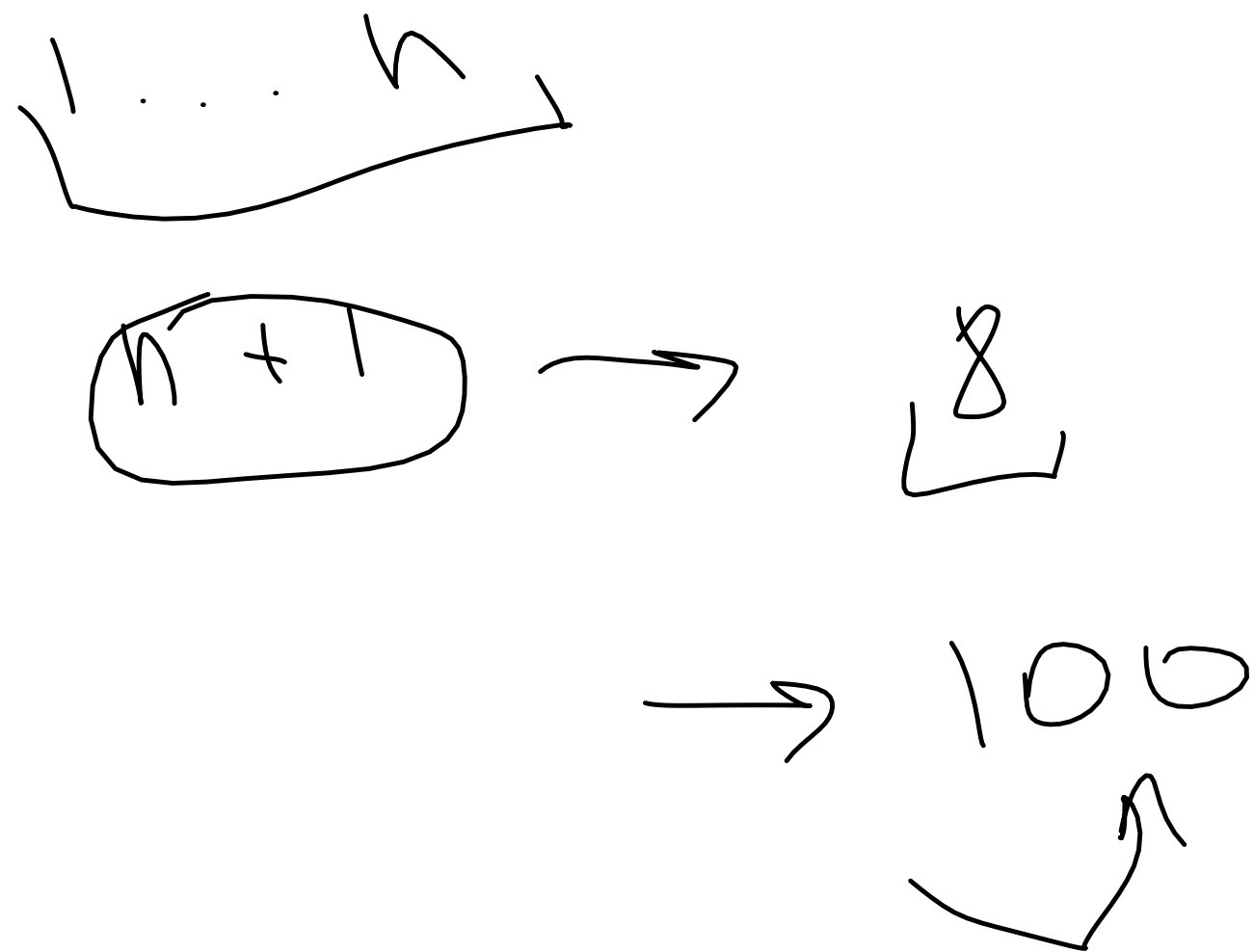
^{key}
Search(item): return item with key if exists

Assumption: items have distinct keys (or that inserting new one clobbers old)

How can we implement the data structure?

Variants:

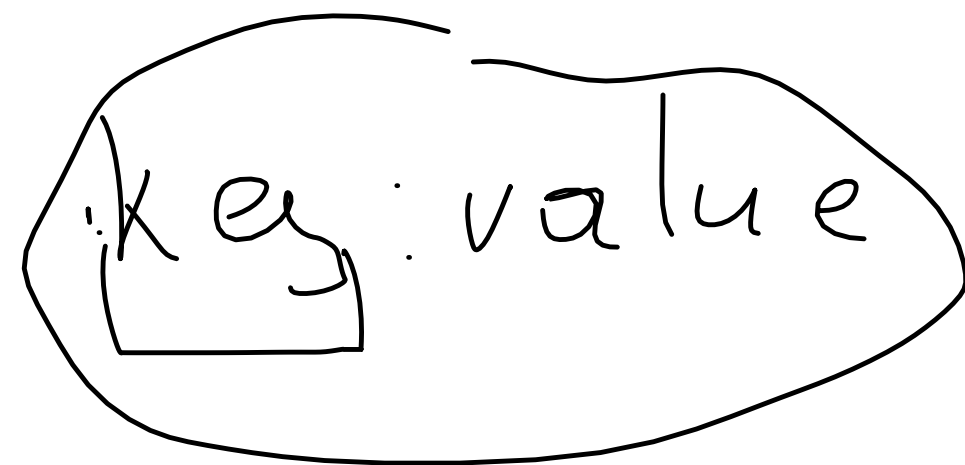
1. *Using arrays.* Solve in $O(n)$ for operation



How can we implement the data structure?

Variants:

1. *Using arrays.* Solve in $O(n)$ for operation
2. *Using Binary Search Trees.* Solve in $O(\log n)$ for operation



How can we implement the data structure?

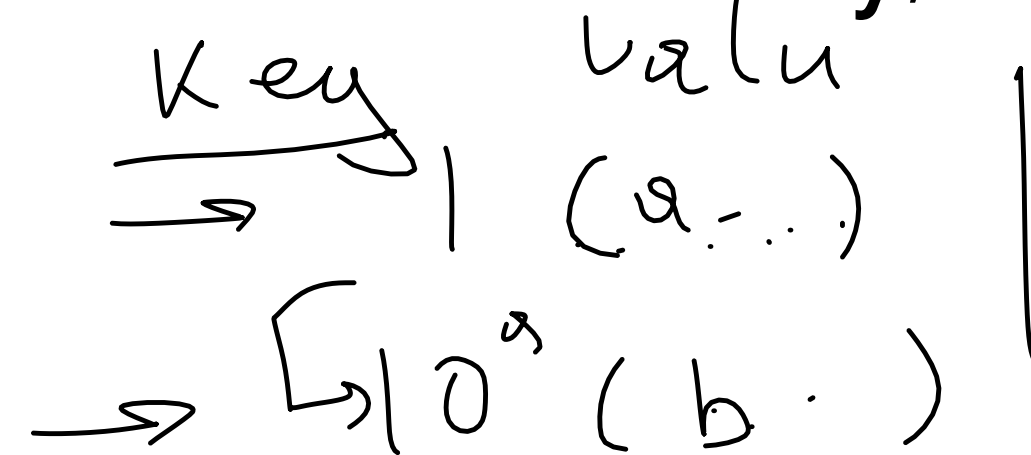
Variants:

1. *Using arrays.* Solve in $O(n)$ for operation
2. *Using Binary Search Trees.* Solve in $O(\log n)$ for operation
3. Our goal is $O(1)$ for operation!

Let's improve the basic approach

Simple approach:

- **Direct access table.** This means items would need to be stored in an array, indexed by key



Problems:

1. Keys must be nonnegative integers (or using two arrays, integers)
2. Large key range \Rightarrow large space e.g. one

key of 2^{256} is bad news

2^{256}

Search(key) Insert(-1, ...)

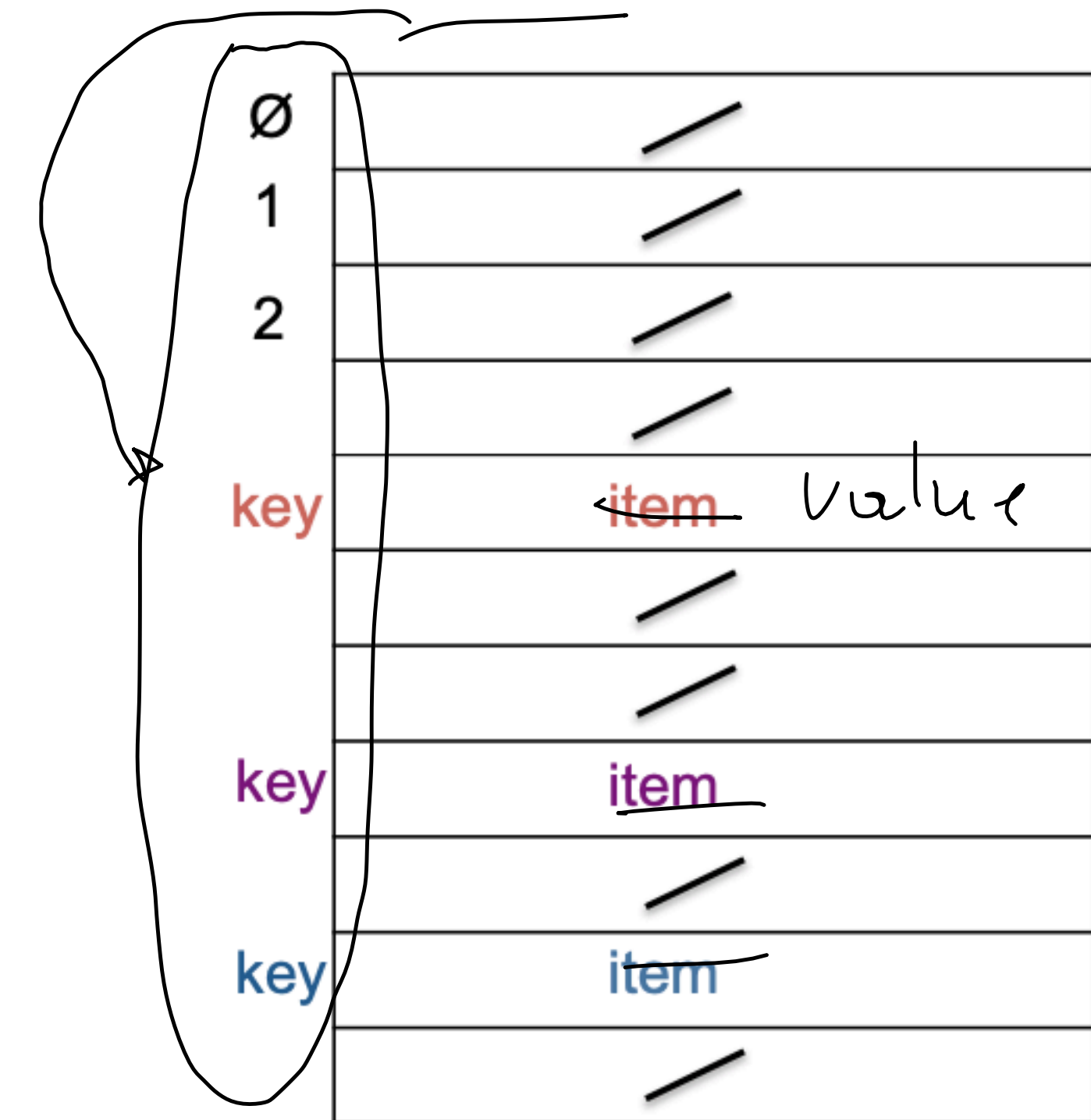


Figure 1: Direct-access table

Let's improve the basic approach

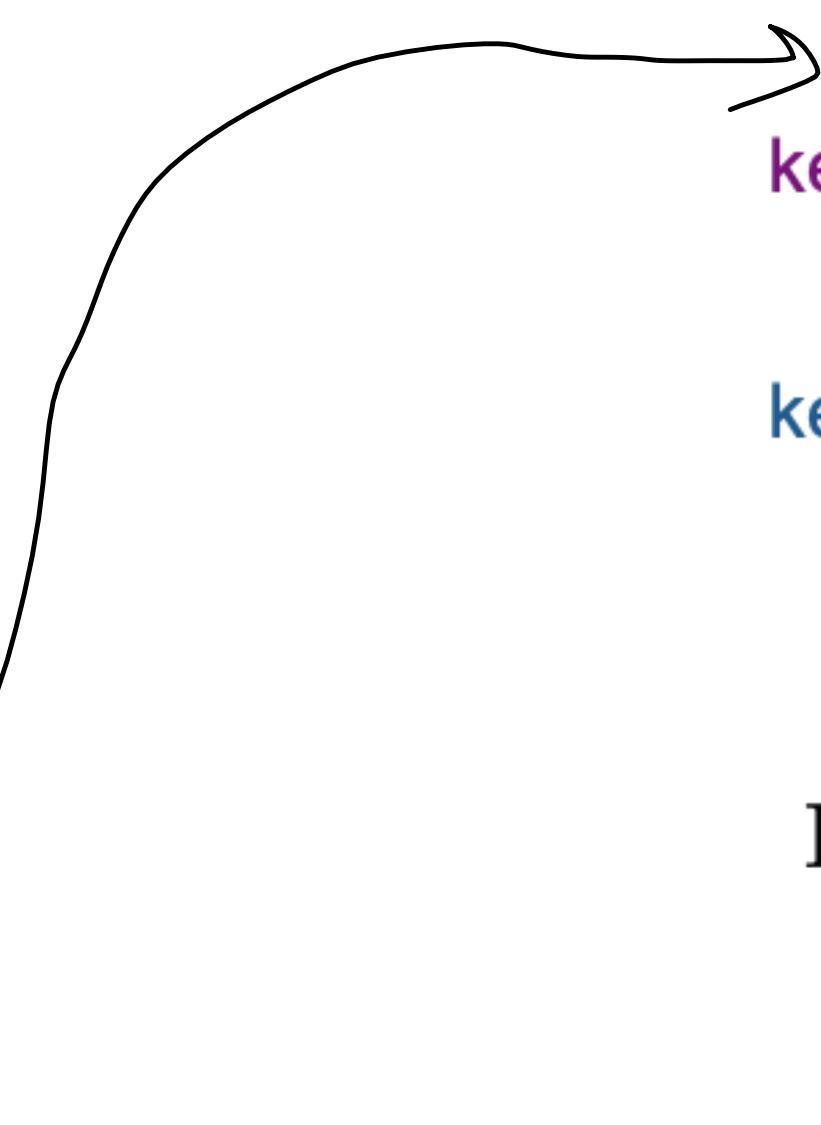
Simple approach:

- **Direct access table.** This means items would need to be stored in an array, indexed by key

Ideas:

1. What if we can map each key to a positive integer number?

$H(\text{John}) \rightarrow (X) \times 70$
 $H(\text{ }) \rightarrow \underline{X, \times 70}$



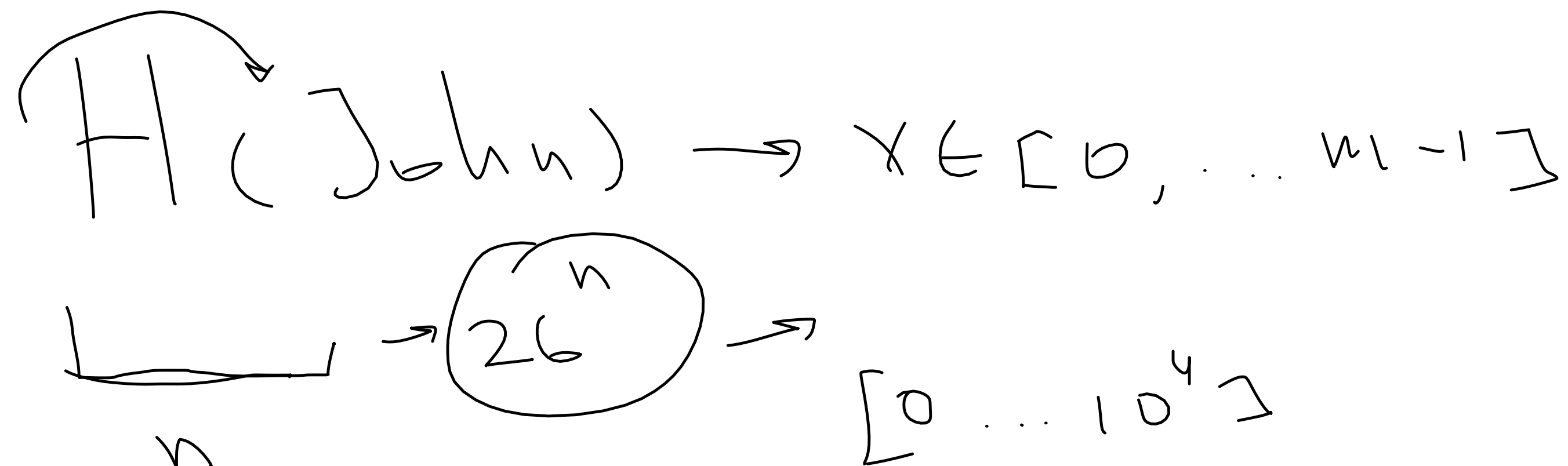
∅	/
1	/
2	/
	/
key	item
	/
	/
key	item
	/
key	item
	/

Figure 1: Direct-access table

Hashing. Hash functions

Definitions:

- **Universe U.** A set of all possible keys
- **Hash function.** $h: U \rightarrow \{0, \dots, m-1\}$
A function that map keys to one of m possible values.
- **Collision**
Keys a,b such that: $a \neq b$ and $h(a) == h(b)$



∅	/
1	/
2	/
	/
key	item
	/
	/
key	item
	/
key	item
	/

Figure 1: Direct-access table

How to use Hashing

Definitions:

- **Universe U.** A set of all possible keys

- **Hash function.** $h: U \rightarrow \{0, \dots, m-1\}$
A function that map keys to one of m possible values.

- **Collision**

Keys a,b such that: $a \neq b$ and $h(a) == h(b)$

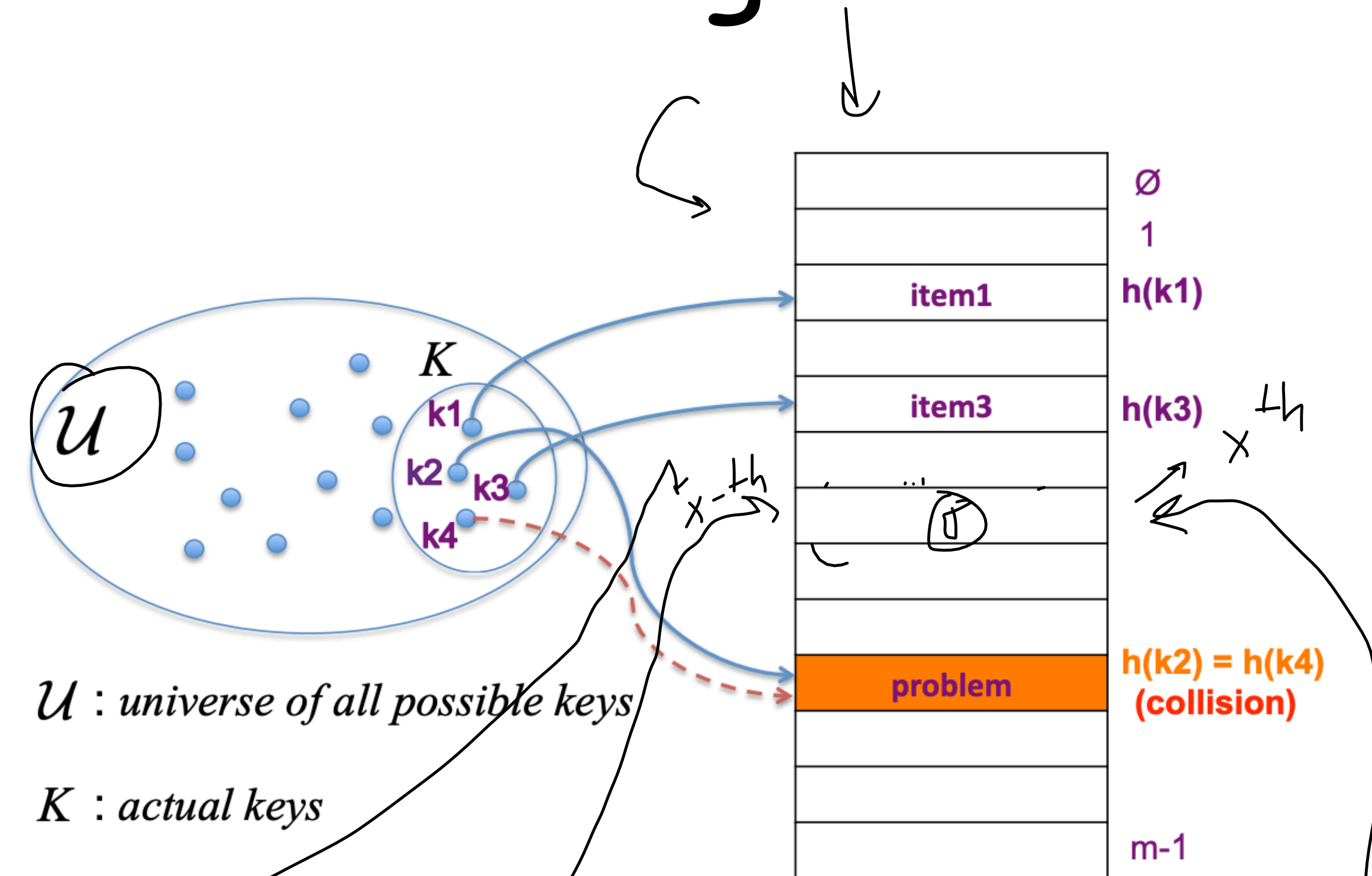


Figure 2: Mapping keys to a table

1) Insert ("aab", 7)
 $h("aab") = x \in [0, \dots, m-1]$

2) Insert ("aba", 7)
 $h("aba") == h("aab") = x$

• two keys $k_i, k_j \in K$ collide if $h(k_i) == h(k_j)$

How to deal with collision?

Chaining:

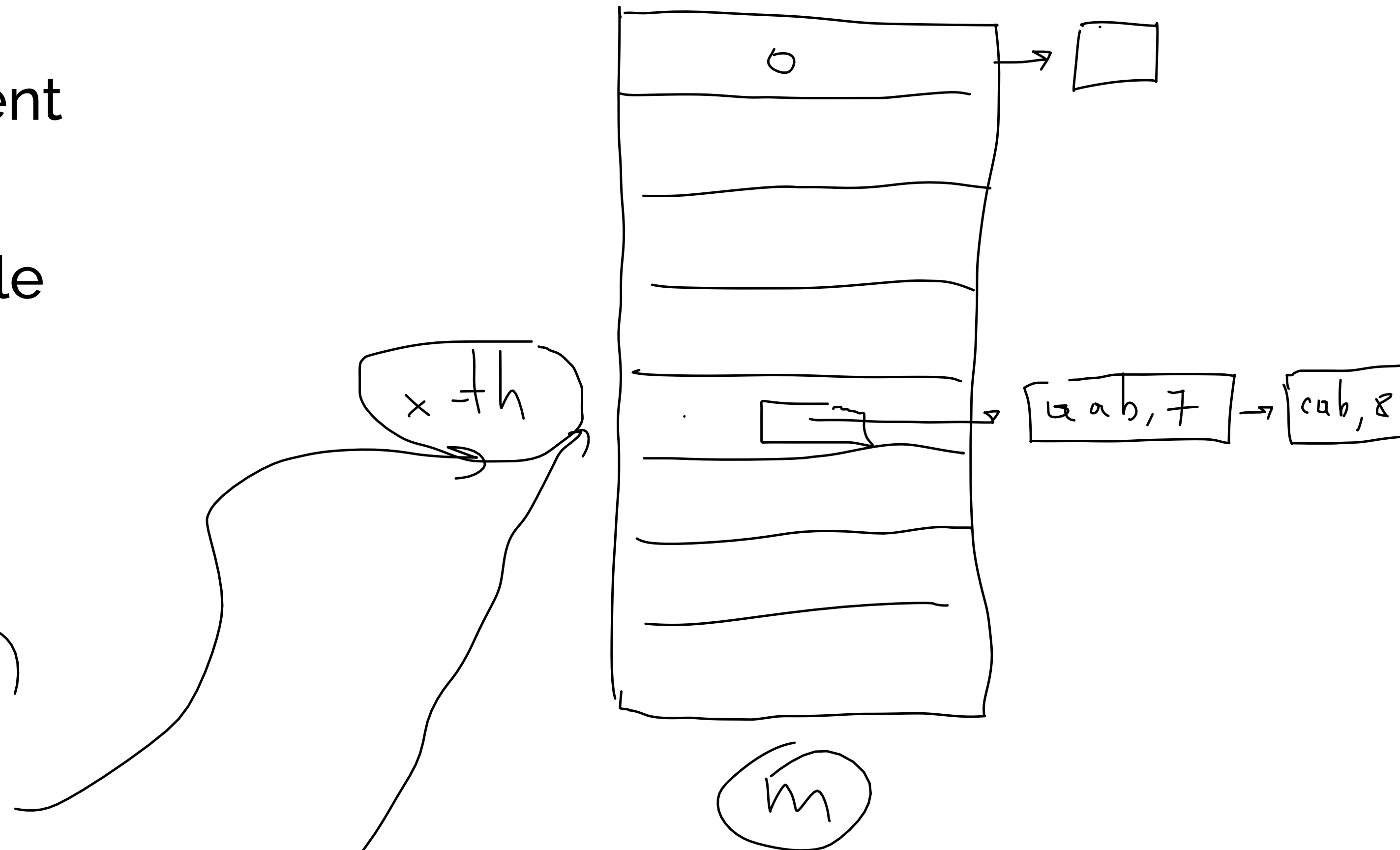
- Linked list of colluding element in each slot of table
- Search must go through whole list $T[h(\text{key})]$

$\text{Insert}(\text{"aab"}, 7)$

$\text{hash}(\text{"aab"}) = x$

$\text{Insert}(\text{"cab"}, 8)$

$\text{hash}(\text{"aab"}) = x$



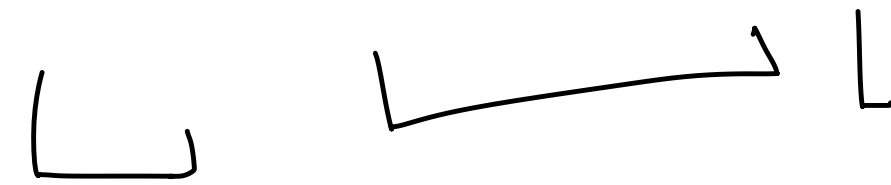
Simple uniform Hashing: Assumption

Load factor:

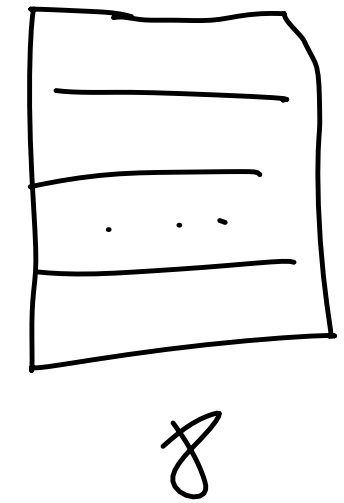
- n - number of keys stored in the table
- m - number of slots in the table
- **Average # keys per slot = n / m = Alpha = load factor**

Expected performance:

- **Search: $O(1 + \text{Alpha})$**
- **Insert/Delete: $O(1 + \text{Alpha})$**
- **The common practice is to keep load factor $\leq 3/4$**
- If load factor $> 3/4 \Rightarrow$ Rehash



$$\frac{20}{8}$$

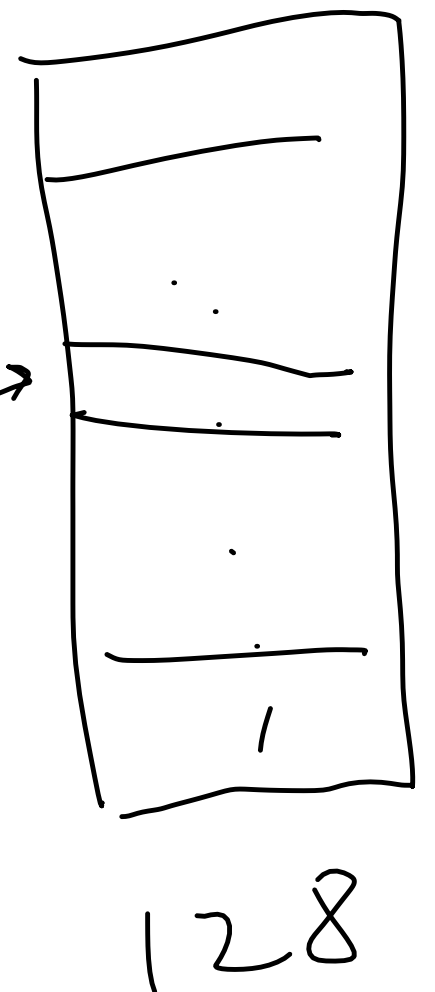


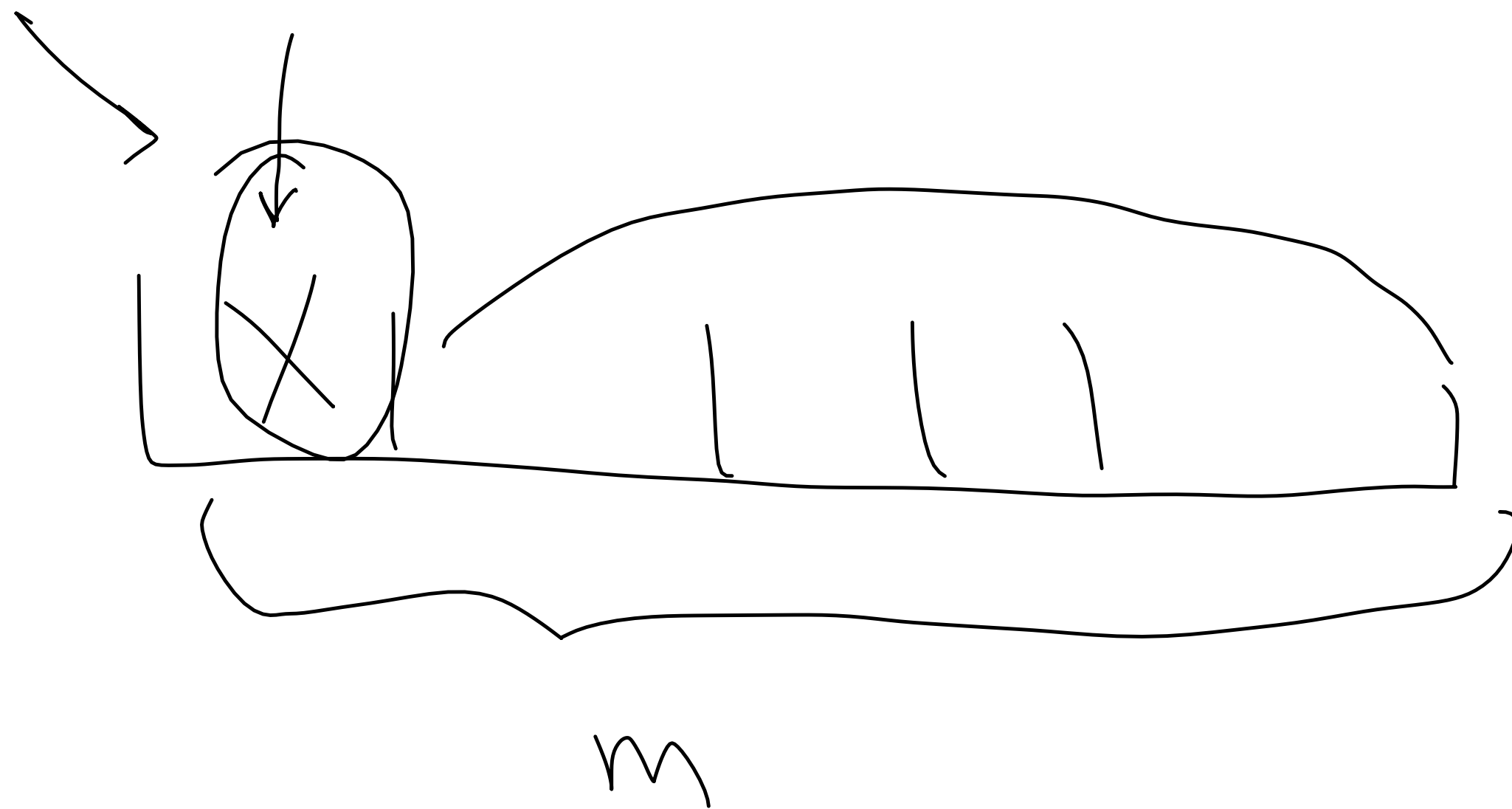
$$n / m = 4$$

hash(..)



$$20/8$$





$$\frac{1}{m-1}$$

Concrete Hash Functions

Division Method: $h(k) = \underline{k \bmod m}$

- k_1 and k_2 collide when $k_1 \equiv k_2 \pmod{m}$, i. e. when m divides $|k_1 - k_2|$
- fine if keys you store are uniform random (probability of collision= $1/m$)
- but if keys are $x, 2x, 3x, \dots$ (regularity) and x & m have common divisor d then use only $1/d$ -th of the table. **Because** $i \cdot x \equiv (i + \frac{m}{d}) \cdot x \pmod{m}$.
(This is likely if m has a small divisor, e. g. 2)
- if $m = 2^r$ then only look at r bits of key! $\boxed{2^r - 1}$
- **Good Practice:** m is a prime number & not close to a power of 2 or 10
(to avoid common regularities in keys)
- **BUT:** Inconvenient to find a prime number; division slow.

Your questions!