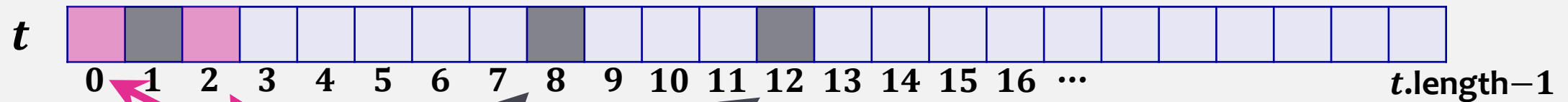# Intro

$$x = 1001110101$$
$$x \bmod 2^7 = 1110101$$

**Hash Tables** are data structures with **very fast** insertion and retrieval data.

almost in constant time



$t$

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15  16  ···  $t.$**length**$-1$

Better to use $t.$**length** that is a prime number not close to a power of **2** or **10**

$data$: **1, 8, 12, 1000000, 18**

$$\mathbf{hash}(x) = x$$
$$\mathbf{hash}(x) = x \bmod \mathbf{10}$$

2

# Intro

id: 9786234291

**Fox**

| | |
|---|---|
| email | fox@forest.ca |
| phone | 613-520-4333 |
| D.O.B | 11.11.1999 |
| address | 8 Blackthorn av. |

with every object there is an associated integer (like id)

# Intro – Open addressing

- **Linear Probing**

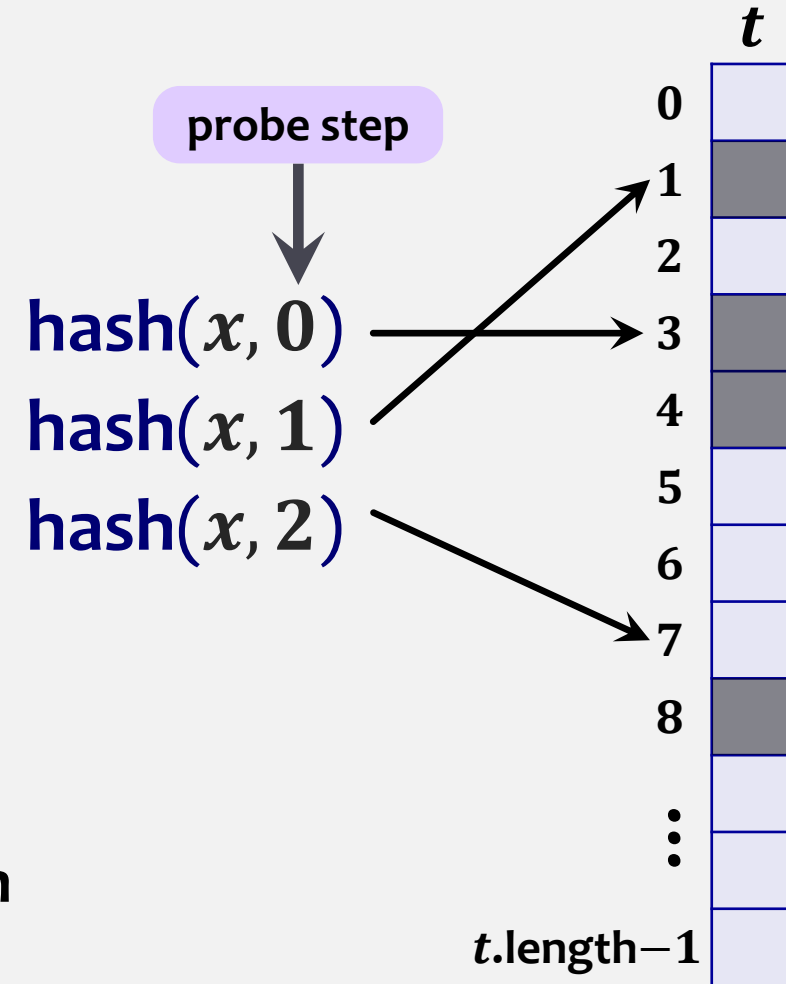  $\text{hash}(x, l) = (\text{hash}(x, 0) + l) \bmod t\text{.length}$

  Problem: **clustering of items**

- **Quadratic Probing**

  $\text{hash}(x, l) = (\text{hash}(x, 0) + l^2) \bmod t\text{.length}$
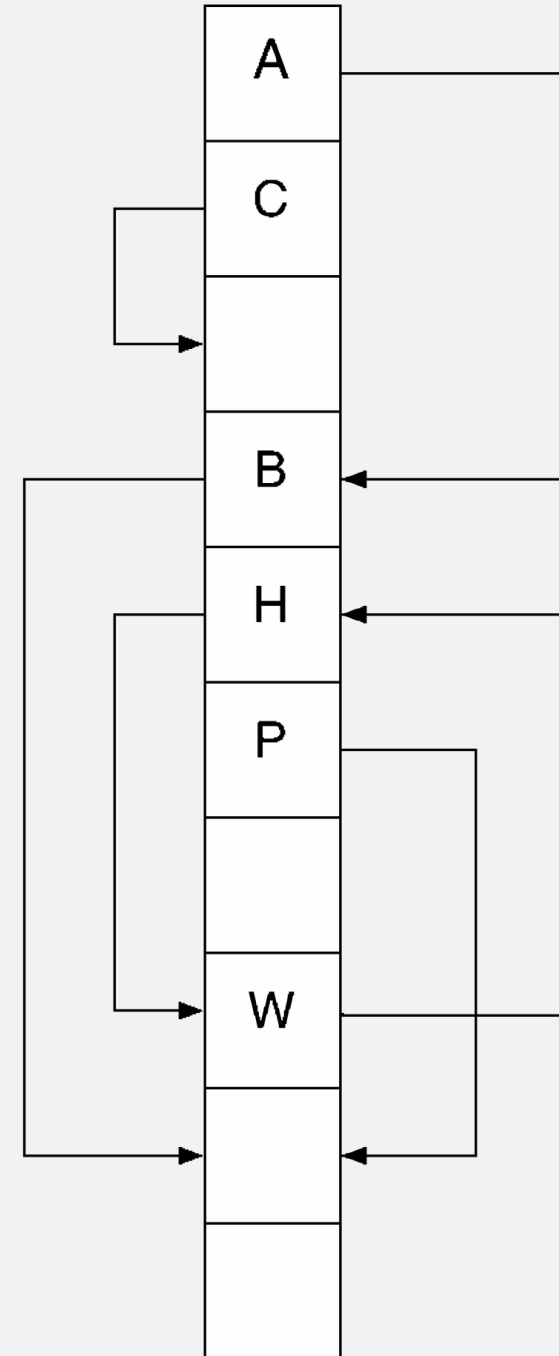
- **Double Hashing**

  $\text{hash}(x, l) = (\text{hash}_1(x) + l \cdot \text{hash}_2(x)) \bmod t\text{.length}$

**Problem:** deleting an item is not very straightforward

$t$

probe step

$\text{hash}(x, \mathbf{0})$

$\text{hash}(x, \mathbf{1})$

$\text{hash}(x, \mathbf{2})$

0
1
2
3
4
5
6
7
8

$t\text{.length}-1$

# Intro – Open addressing

- **Cuckoo hashing**

# Hash Tables

- One of the most used data structuring techniques.

- Largely misunderstood (the most problematic area is computer security)

- Implement **Set** (and **Map**) interface:

  - add($x$)
  - remove($x$)
  - contains($x$)

$$O(1) \text{ expected time}$$
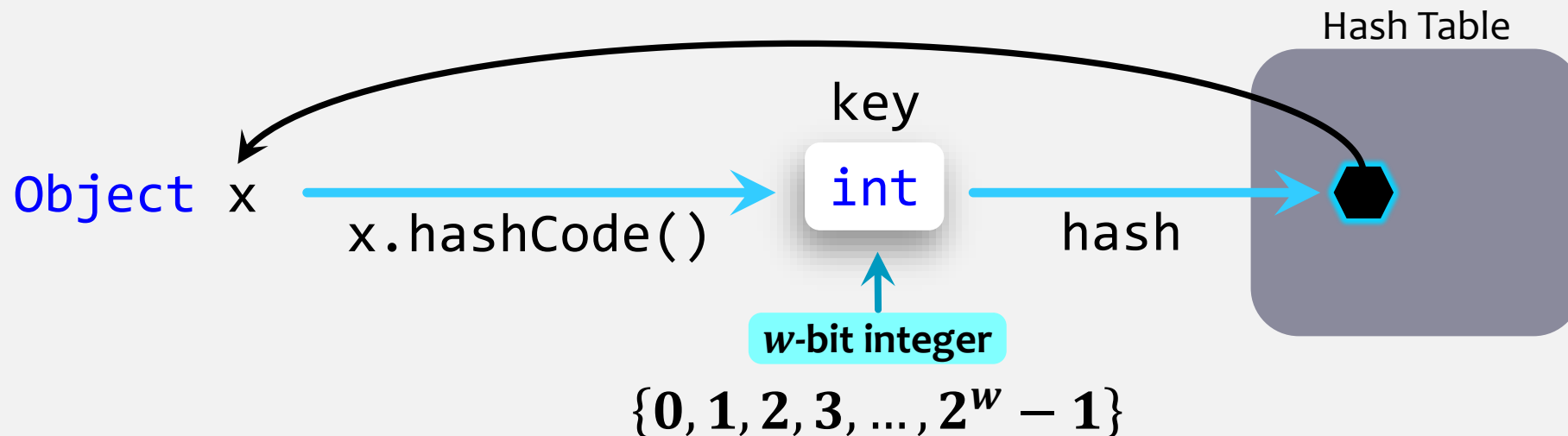
**independent of the size of the set**

- Make some internal **random** choices.

# Hash Tables

- Hash tables are an efficient method of storing a small number of **integers** from a large range.

- Very often hash tables store types of data that are **not integers**:

  1. Given an object, you associate it with an integer that is suitable for storing in a hash table. In Java:  `x.hashCode()`

  2. Hash Table will store that integer (together with a reference to $x$)

**hashCode() method (class Object) returns a hash code of the object  x**

Hash Table

key

Object x

x.hashCode()
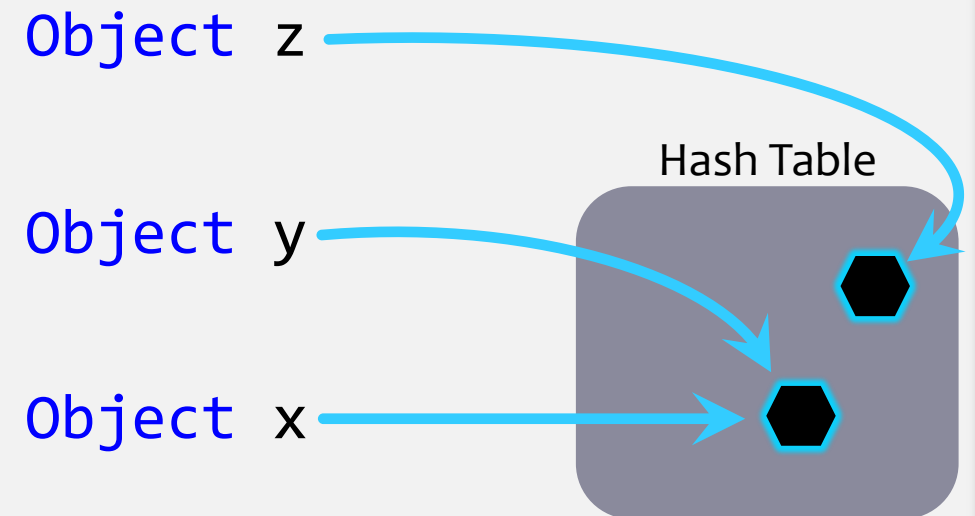
int

hash

*w*-bit integer

$\{0, 1, 2, 3, \ldots, 2^w - 1\}$

# Common Mistake

If two or more objects are equal according to the **equals** method (Object class), then their hashes should be equal too:

$$\texttt{if x.equals(y) then x.hashCode() == y.hashCode()}$$

If you override the **equals** method, it is crucial
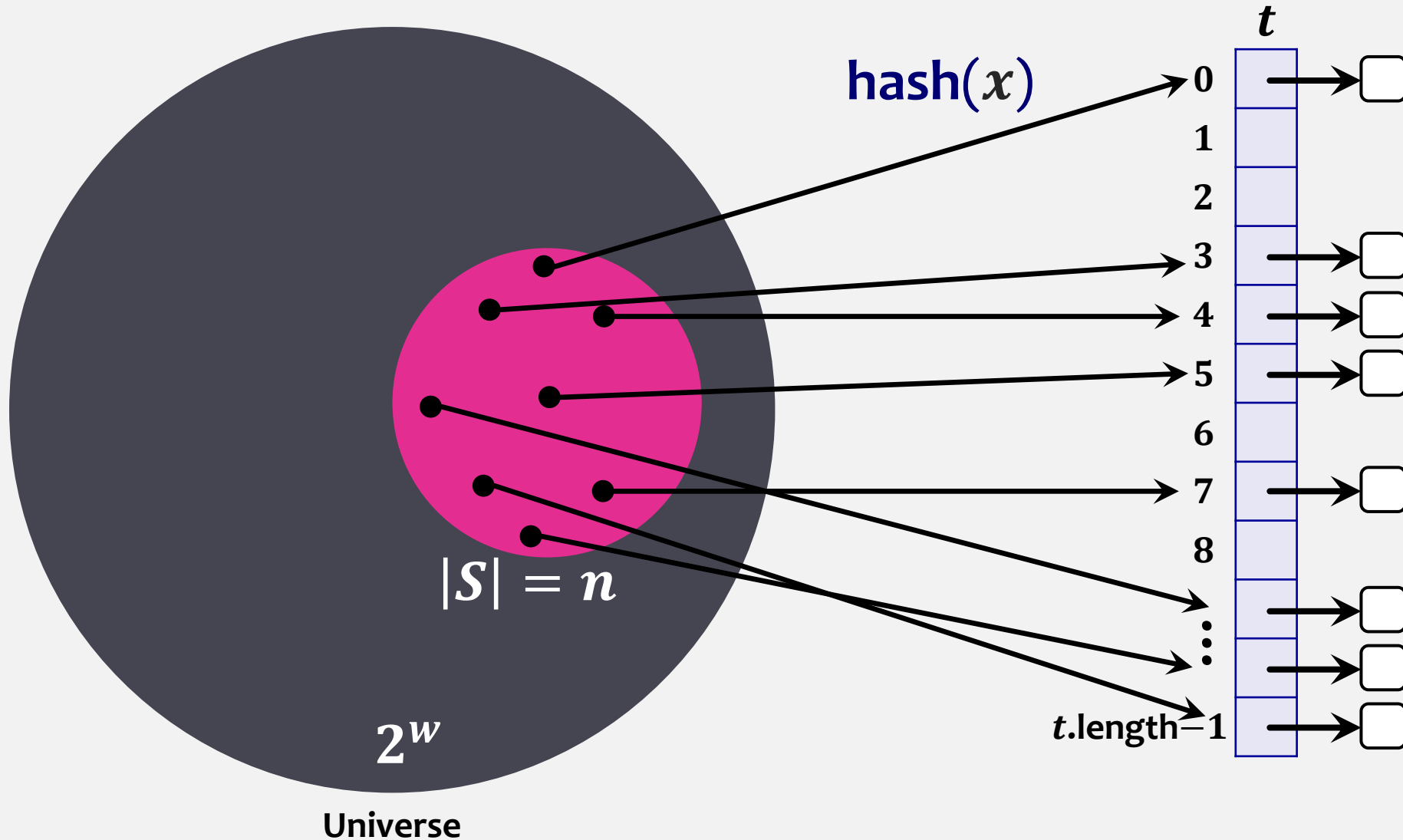to override the **hashCode** method as well.

Note: If two objects are **not equal**, we want their
hashes to be **not equal**.
But in reality, they can be **equal or unequal**.

`Object z`

`Object y`

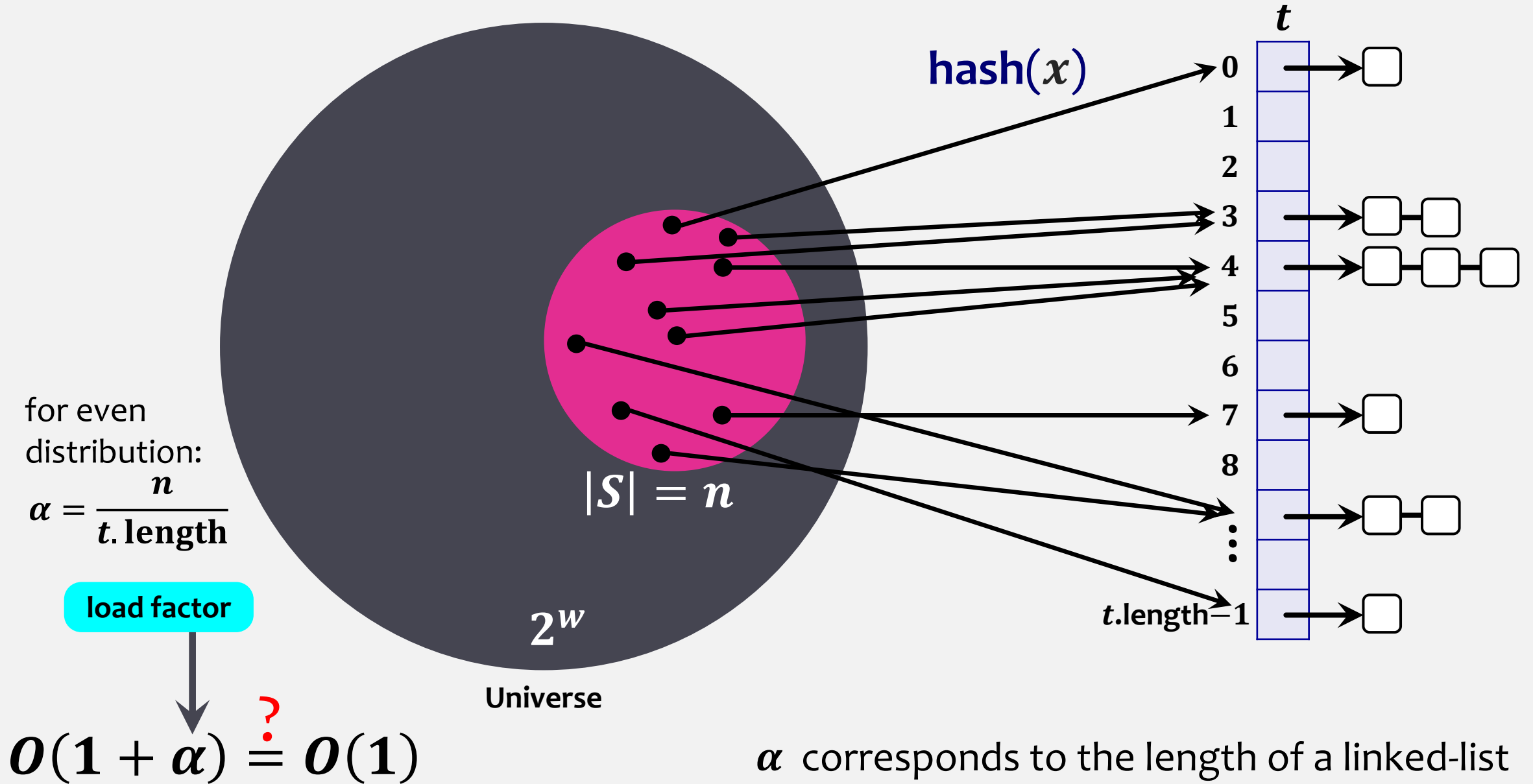`Object x`

Hash Table

Passwords.java & Hashes.java

https://www.educative.io/edpresso/what-is-the-hashcode-method-in-java

# Chained Hash Table

**hash**$(x)$

$t$

0
1
2
3
4
5
6
7
8

$t$.length$-1$

$|S| = n$

$2^w$

**Universe**

$O(1)$     **Ideal situation**

# Chained Hash Table



**hash**$(x)$

$t$

for even distribution:
$$\alpha = \frac{n}{t.\text{length}}$$

**load factor**

$$O(1 + \alpha) \overset{?}{=} O(1)$$

$|S| = n$

$2^w$

**Universe**

$t.$length$-1$

$\alpha$ corresponds to the length of a linked-list
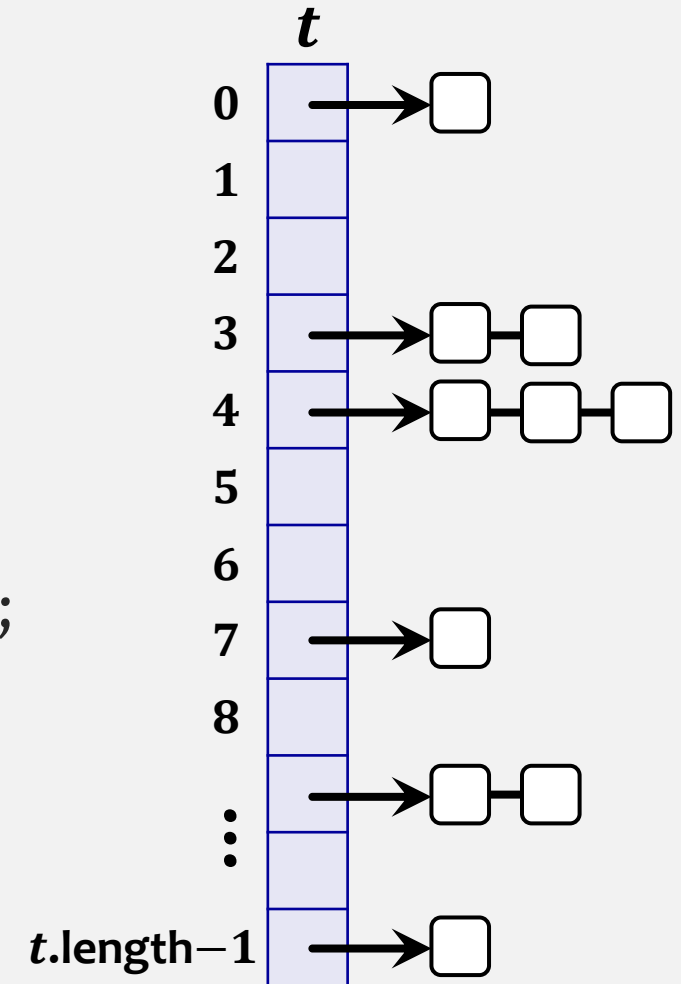
# Chained Hash Table

- Assume: the set of objects we're trying to store is a **set of integers** that are all **distinct**.

- Hash table is an array $t$ of lists.

- $n$ is the total number of items in all lists.

- All items with hash value $i$ are stored in the list at $t[i]$.

- function $\mathbf{hash}(x)$ returns the hash value of a data item $x$; $x$ is a $w$-bit integer (for now): $x \in \{0, 1, \ldots, 2^w - 1\}$; hash value $i$ is in the range $\{0, 1, 2, \ldots, t.\mathbf{length} - 1\}$.

- for lists not to get too long, we maintain $n \leq t.\mathbf{length}$

- average number of elements stored in one of these lists is 
$$\frac{n}{t.\mathbf{length}} \leq 1$$

# contains($x$)

We perform a linear search on the list $t[\text{hash}(x)]$:

T contains($x$):

    $i = \textbf{hash}(x)$;

    $\textbf{list} = t[i]$;
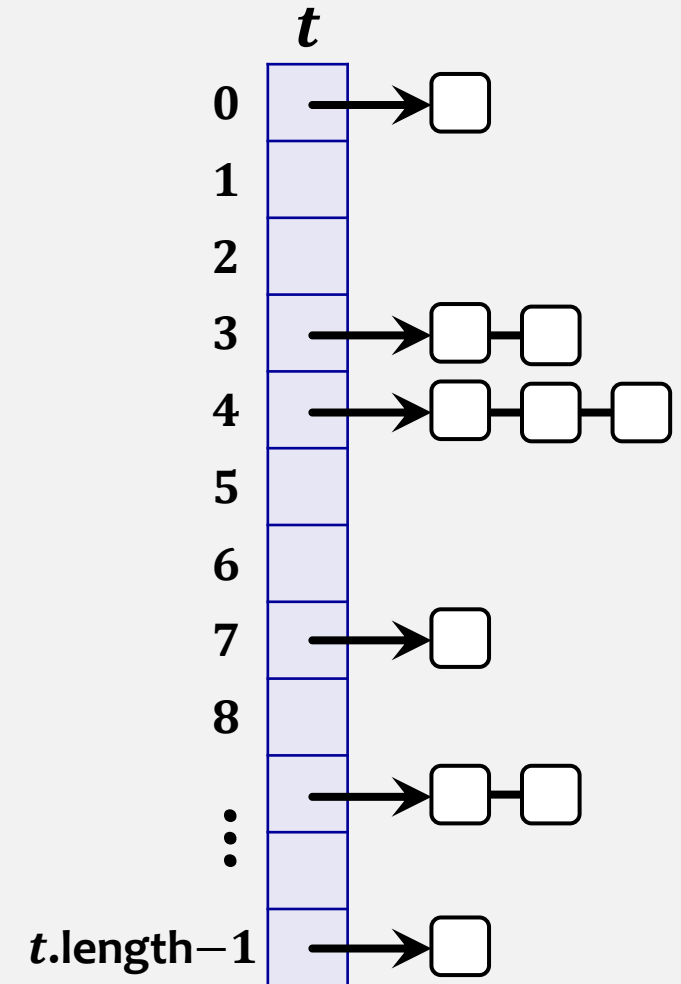
    for each ($y$ in **list**)

        if ($y$.equals($x$)) then

            return $y$;

    return null;

$$O\left(n_{\text{hash}(x)}\right)$$ Where $n_{\text{hash}(x)}$ is the size of the list $t[\textbf{hash}(x)]$

# add($x$)

the cost of growing is only constant when amortized over a sequence of insertions

$O(n)$ • If the length of $t$ needs to be increased, then grow $t$.

$O(1)$ • hash $x$ to get an integer $i \in \{0, 1, 2, \dots, t.\text{length} - 1\}$

$O(1)$ • <mark>append</mark> $x$ to the list $t[i]$

**remove this if you want to store the same $x$ again**

boolean add($x$):

    if (contains($x$) $\neq$ null) then return false;

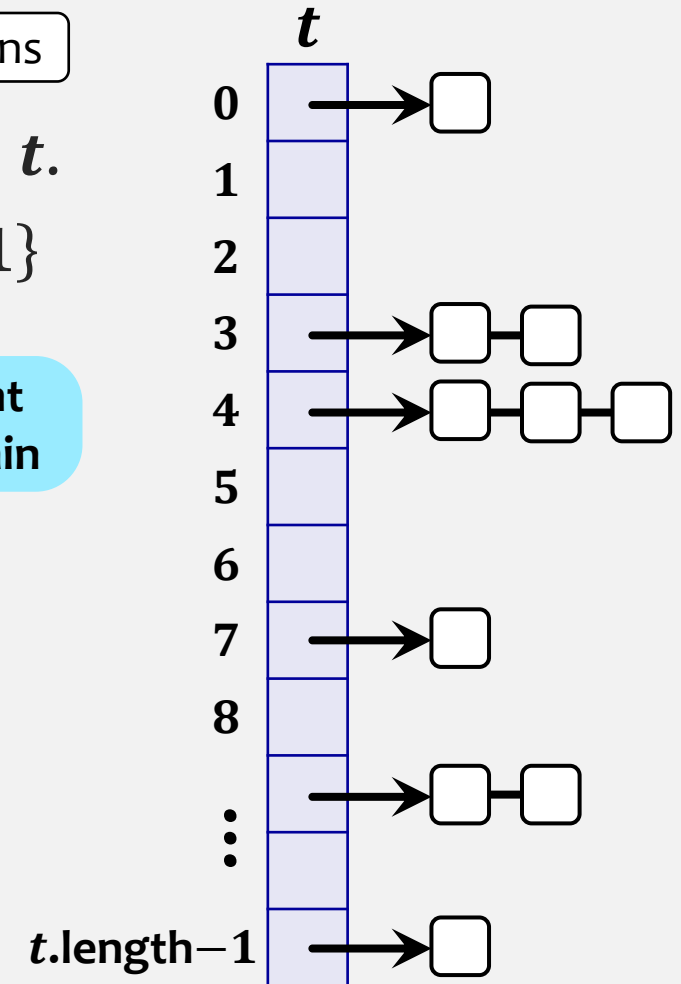    if ($n + 1 > t.\text{length}$) then resize();

    $t[\textbf{hash}(x)]$.add($x$);

    $n + +$;

    return true;

$i = \textbf{hash}(x)$;
$\textbf{list} = t[i]$;
$\textbf{list}.\text{add}(x)$;

$O(n_{\text{hash}(x)})$



t

0
1
2
3
4
5
6
7
8
$t.\text{length}-1$

# remove($x$)

- hash $x$ to get an integer $i \in \{0, 1, 2, \ldots, t.\text{length} - 1\}$
- iterate over the list $t[i]$ until you find $x$, and remove it
- (optional) if the length of $t$ needs to be decreased, then shrink $t$.



T remove($x$):

    if (contains($x$) = null) then return null;

**expensive** →    $y = t[\text{hash}(x)]$.remove($x$);

    $n--$;
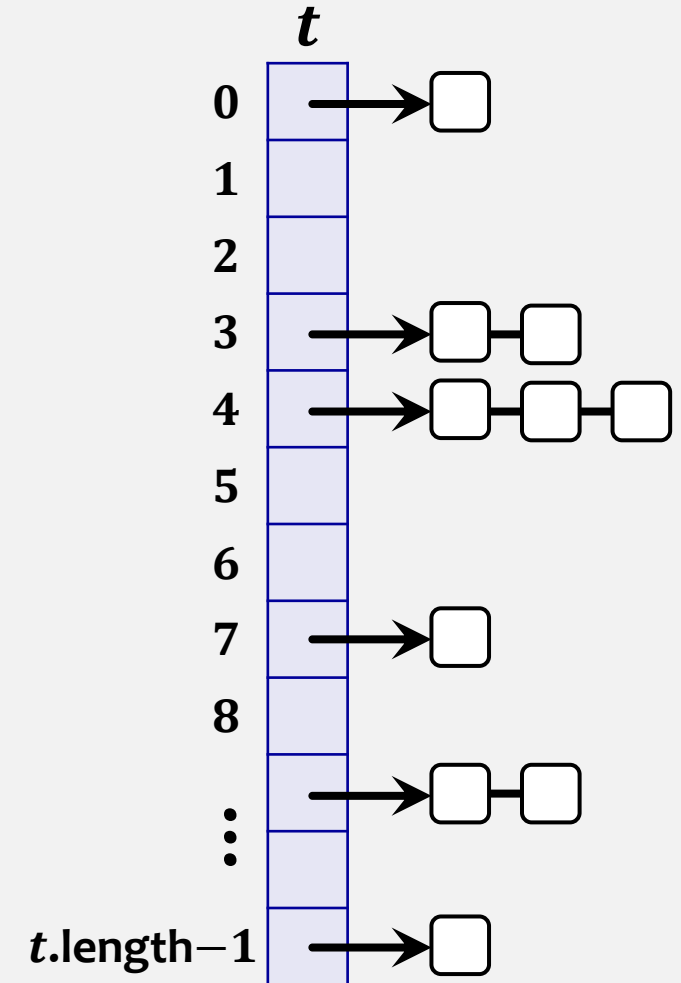
    if ($4n < t.\text{length}$) then resize();

    return $y$;

$O(n_{\text{hash}(x)})$    Where $n_{\text{hash}(x)}$ is the size of the list $t[\text{hash}(x)]$

16

# Hash Functions – Good & Bad
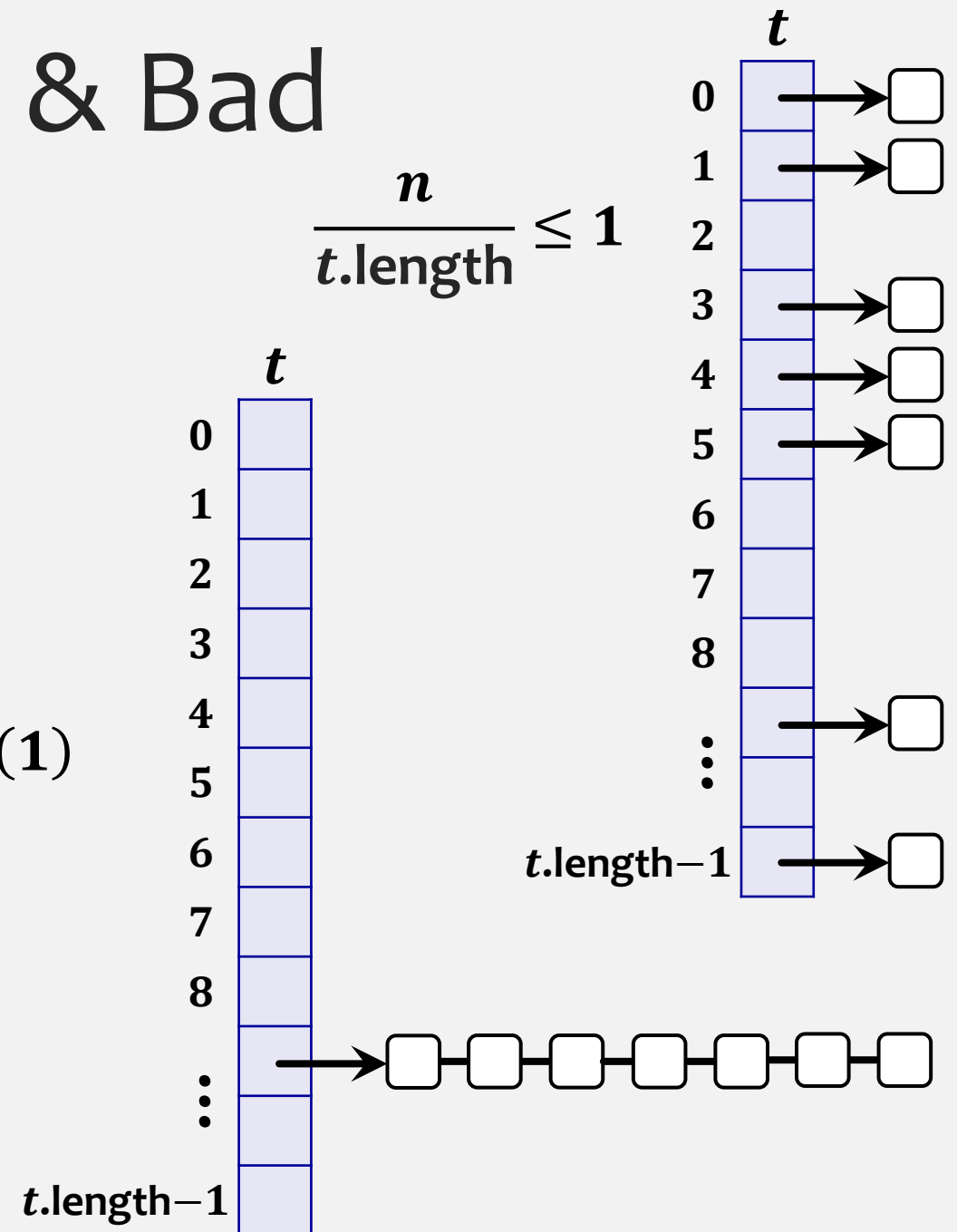
The performance of a hash table depends on the choice of the hash function.

$$\frac{n}{t.\textbf{length}} \leq 1$$

- A **good** hash function will spread the elements evenly among the $t.\textbf{length}$ lists, so the expected size of the list $t[\textbf{hash}(x)]$ is

$$O\left(\frac{n}{t.\textbf{length}}\right) = O(1)$$

  **good** hash function should not depend on patterns in the data

- A **bad** hash function will hash all values (including $x$) to the same table location, so the size of the list $t[\textbf{hash}(x)]$ will be $n$.

17

# Universal Hashing

We select a hash function **at random** from a family of hash functions with a certain mathematical property. This guarantees a low number of collisions in expectation.

We want for any $x, y \in \{0, \dots, 2^w - 1\}, \; x \neq y$ $\qquad Pr(\mathbf{hash}(x) = \mathbf{hash}(y)) \leq \dfrac{1}{t.\mathbf{length}}$

Example of a family of hash functions:

Let $z$ be a random number in $\{0, \dots, 2^w - 1\}$, and let $t.\mathbf{length}$ be a prime number.

Then the formula for hashing an integer $x$ is

Carter & Wegman - 1979 $\longrightarrow$ $\mathbf{hash}(x) = (z \cdot x) \bmod t.\mathbf{length}$

expensive

# Family of hash functions

$$\text{hash}(x) = (z \cdot x) \bmod t.\text{length}$$

**Improvement:**  let $t.\text{length}$ be $2^k$ (for some integer $k$)

$$y \bmod 2^k \quad \equiv \quad \text{last } k \text{ bits of } y$$

**Problem:**  there are certain sets of integers that are hashed with lots of collisions. With this hash function you can get lists of logarithmic length.