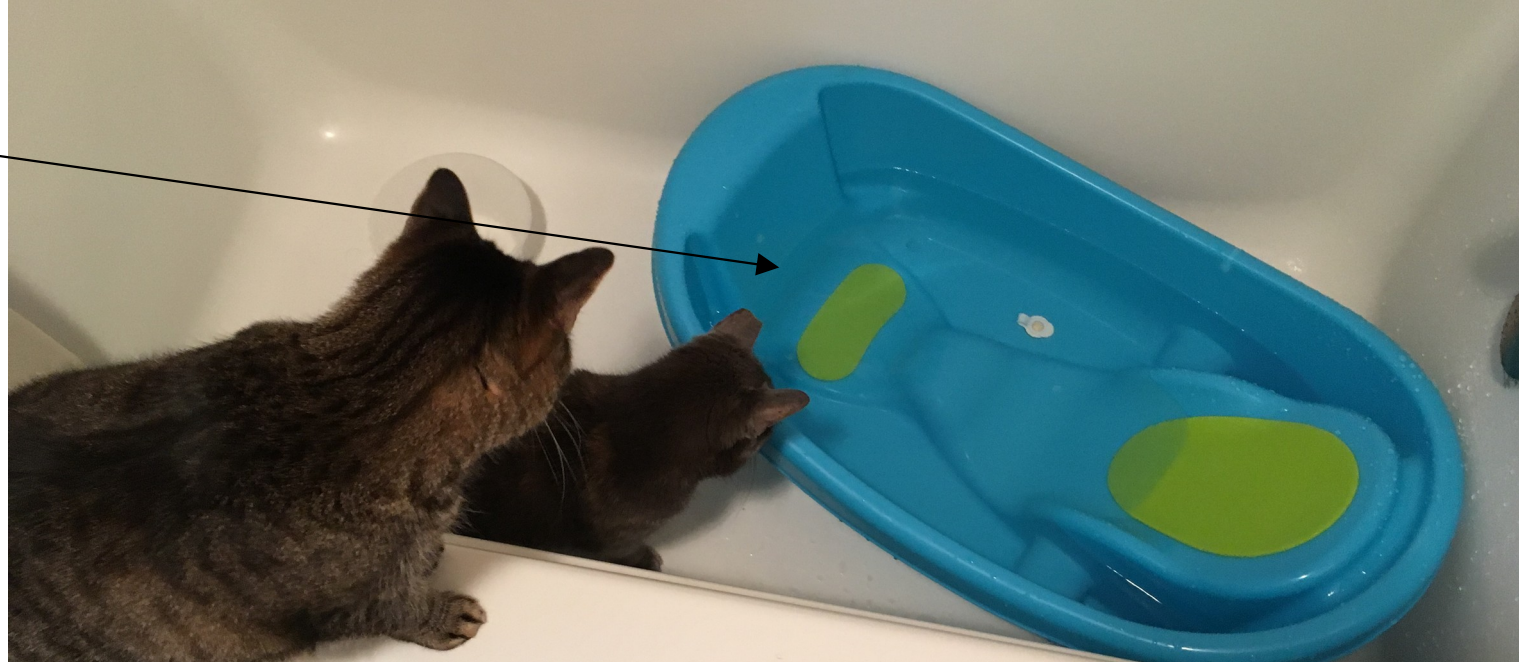An unacceptable position



# CSE 250
# Lecture 29

Hash Tables

# Maps

- Sets, Bags: Collections of Elements of type A
    - add(a: A)
    - remove(a: A)
    - apply(a: A)
        - set.apply(a:A): Boolean // is the element part of set
        - bag.apply(a:A): Int // # of copies of the element in bag
- Maps: Like Sets, but where A is a 2-tuple (key, value)
    - The identity of the element is determined by key

# Maps

- Map[K, V]
  - add(k: K, v: V)                 // also called put(k, v)
    - Insert (k, v) into the map
    - If an element associated with key k already exists, replace it.
  - remove(k): V
    - Remove the element associated with key k, return the corresponding value
  - apply(k: K): V                 // also called get(k)
    - Return the value corresponding to key k.

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Maps

- Map[K,V] as a Sorted Sequence

    - O(log(n)) apply (but very cache-friendly)

    - O(n) add

    - O(n) remove

- Map[K,V] as a balanced Binary Search Tree

    - O(log(n)) apply

    - O(log(n)) add
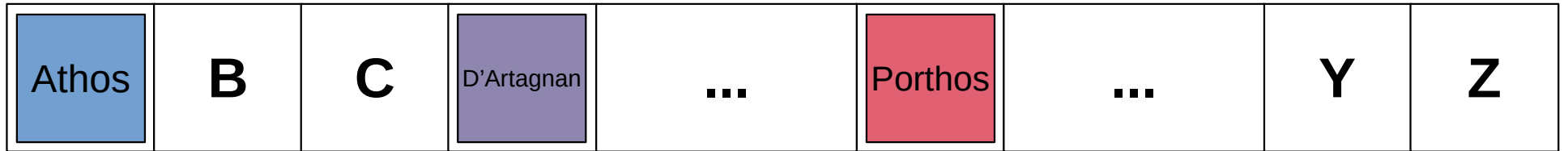
    - O(log(n)) remove

- Map[K,V] as a LSM tree

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Finding Items

The most epensive part of finding records is **finding** them.
(i.e., where is the record located?)

**So... skip the search**

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Alternative Idea: Assign Bins

- Create an array of size N

- Pick an O(1) function to assign each record a number in [0,N)
  - First letter of name → [0, 26)

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Alternative Idea: Assign Bins

| Athos | B | C | D'Artagnan | ... | Porthos | ... | Y | Z |

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Alternative Idea: Assign ~~Bins~~ Buckets

- Pros

  - O(1) Insert

  - O(1) Find

  - O(1) Remove

- Cons

  - Wasted Space (Only 3/26 slots used)

  - Duplication (What about Aramis?)

©Oliver Kennedy, Eric Mikida, Andrew Hughes
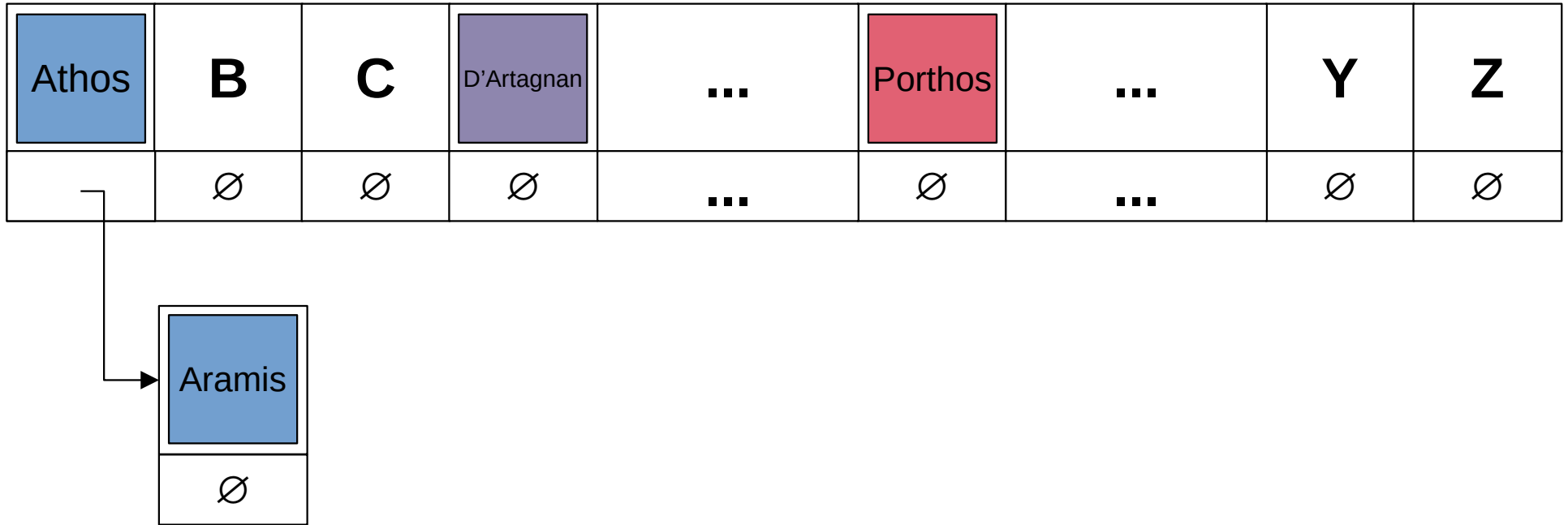The University at Buffalo, SUNY

# Bucket-Based Organization

- Wasted Space
  - Not ideal, but not wrong
  - O(1) access time might be worth it!
  - Also depends on choice of function (more on this later)
- Duplication
  - We need to deal with duplicates!

# Duplication

- **Idea**: Make buckets bigger (e.g., B elements)

  – **Pro**: Up to B duplicates in a bucket, Still O(1) (O(B)) find

  – **Con**: No more than B duplicates in a bucket

- **Idea**: Make buckets arbitrarily large (e.g., Linked List)

  – **Pro**: No more overflow

  – **Con**: O(n) worst-case find.

# Buckets + Linked Lists

| Athos | **B** | **C** | D'Artagnan | **...** | Porthos | **...** | **Y** | **Z** |
|---|---|---|---|---|---|---|---|---|
| | ∅ | ∅ | ∅ | ... | ∅ | ... | ∅ | ∅ |

Aramis

∅

©Oliver Kennedy, Eric Mikida, Andrew Hughes
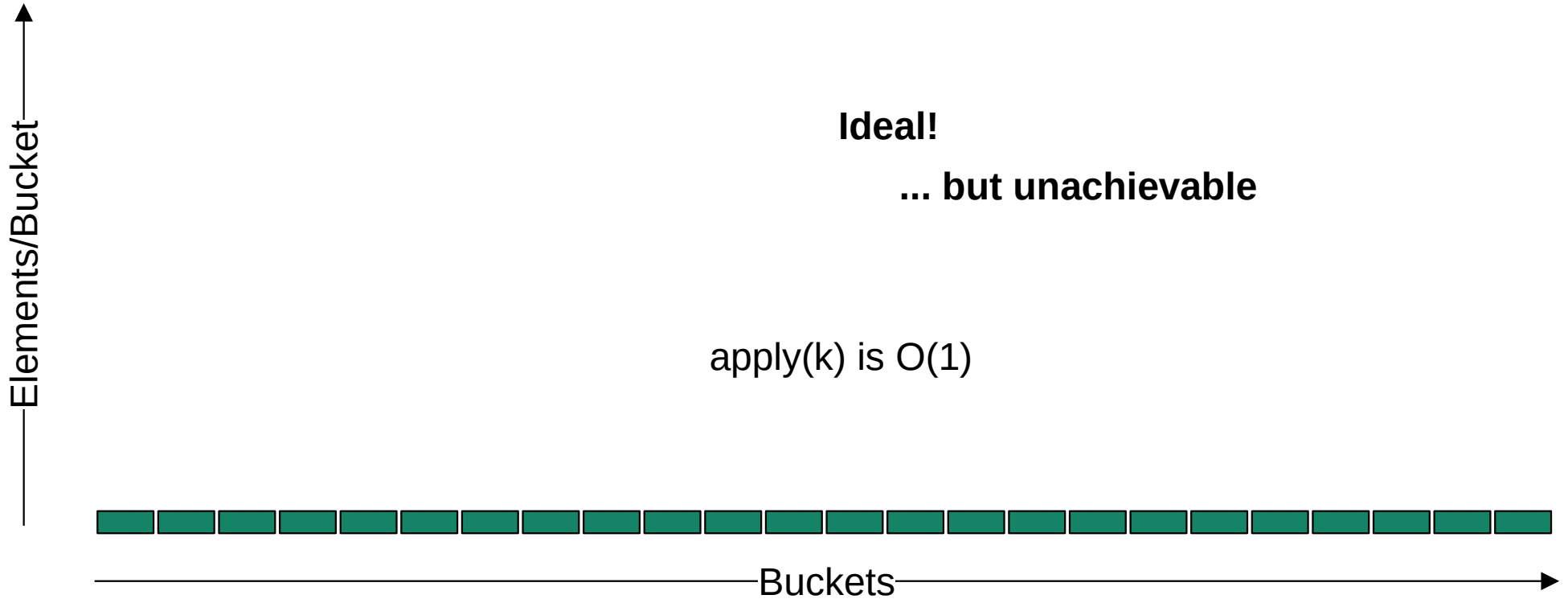The University at Buffalo, SUNY

# Buckets + Linked Lists

```scala
class BucketMap[K, V](_numBuckets: Int, _lookupFunction: K => Int) {

  val _buckets = new Array[ List[(K, V )] ](_numBuckets)
  for(i <- 1 until _numBuckets) { _buckets(i) = Nil }

  def apply(key: K): V = {
    val bucketPosition = _lookupFunction(k)
    var element = _buckets(bucketPosition)
    while(element != Nil){
      if(key == element.head._1) { return element.head._2; }
      element = element.tail
    }
    throw new NoSuchElementException(key.toString)
  }
}
```
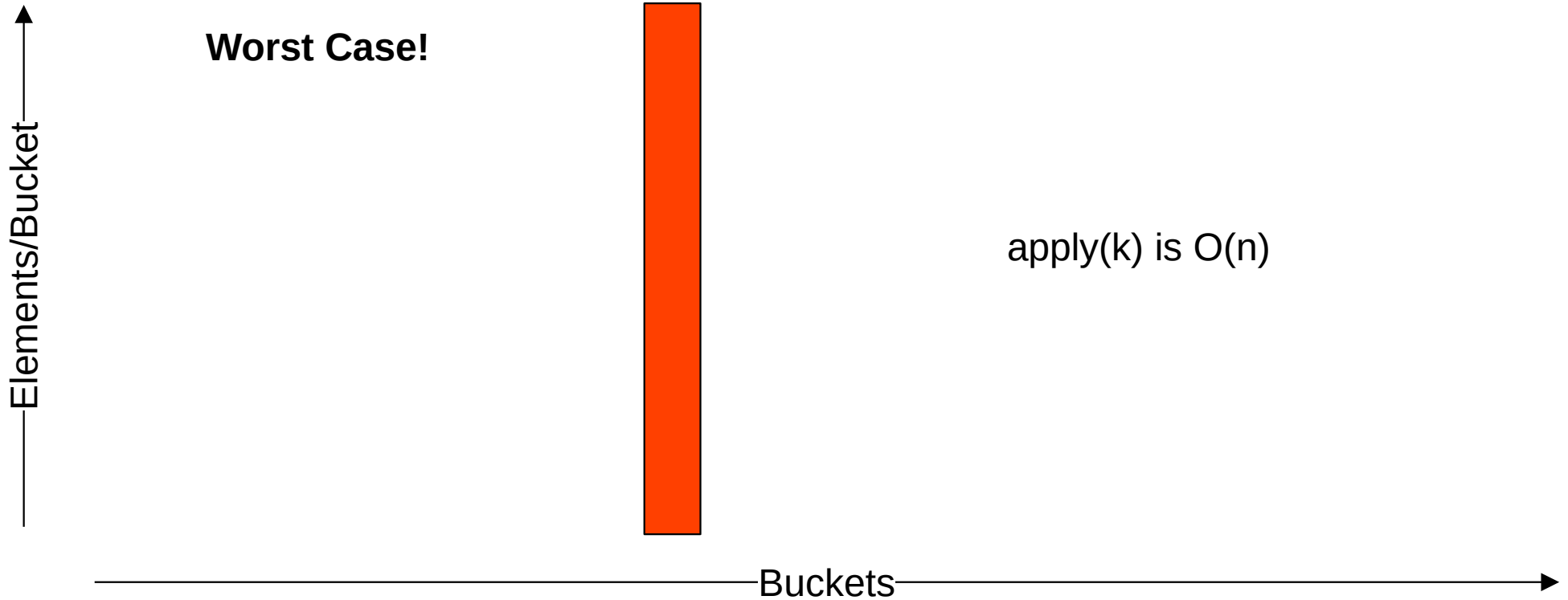
©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Picking a Lookup Function

- Desirable Features for h($x$)

  - Fast

    - needs to be O(1)

  - "Unique"

    - As few duplicate bins as possible

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Picking a Lookup Function

Elements/Bucket →

**Ideal!**

**... but unachievable**

apply(k) is O(1)

Buckets →

# Picking a Lookup Function

**Worst Case!**

Elements/Bucket

apply(k) is O(n)

Buckets

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Picking a Lookup Function

**Almost Ideal!**

**... and achievable**

apply(k) is something like O(1)?

Elements/Bucket

Buckets

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Other Functions

- First Letter of UBIT

# First letter of UBIT



36 'j's

No 'u's

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Other Functions

- First Letter of UBIT name

  - Unevenly Distributed: O(n) apply

- Identity Function on UBIT #

  - Need a 50m+ element array

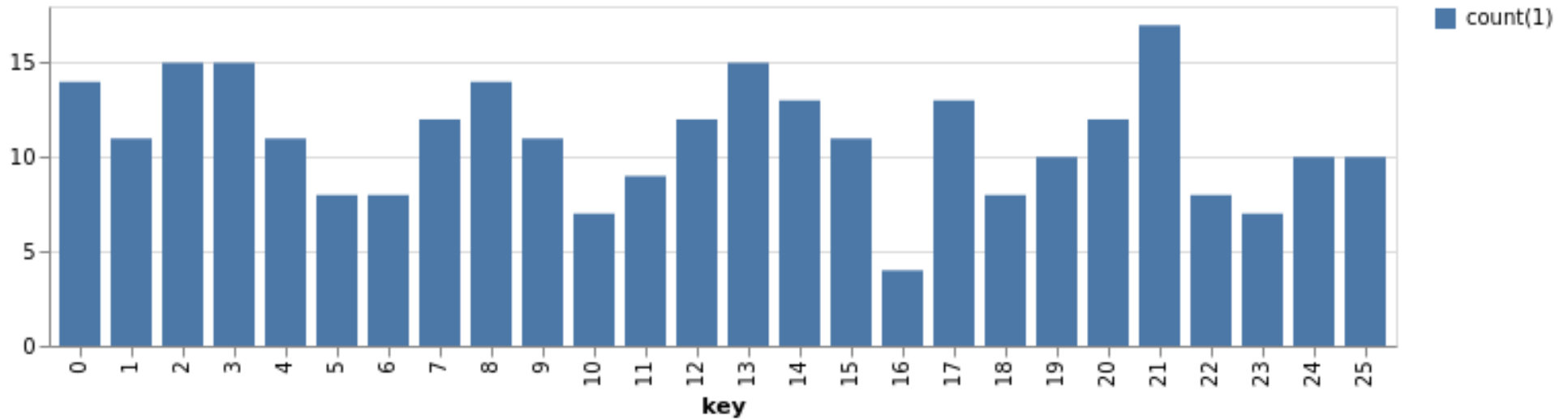©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Other Functions

- Identity Function: `(x: Int) => x`

  - **Problem**: Can return values over N

  - **Solution**: Cap return value by Modulus with N

    - `(x: Int) => x % N`

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Modulus

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| 14 | 15 | 16 | 17 | 18 | 19 | 20 |

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Modulus

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Modulus

**But still relies on UBIT # being random!**

UBIT # % 26

substr(UBITName, 0, 1)

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Picking a Lookup Function

- **Wacky Idea**: Have h(x) return a random value in [0, N)

  - Random.nextInt % N

    (Yes, it makes apply impossible, but bear with me)

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Picking a Lookup Function

$$n = \text{number of elements in any bucket}$$

$$N = \text{number of buckets}$$

$$b_{i,j} = \begin{cases} 1 & \textbf{if } \text{element } i \text{ is assigned to bucket } j \\ 0 & otherwise \end{cases}$$

$$\mathbb{E}\left[b_{i,j}\right] = \frac{1}{N}$$

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Picking a Lookup Function

$$n = \text{number of elements in any bucket}$$

$$N = \text{number of buckets}$$

$$b_{i,j} = \begin{cases} 1 & \textbf{if } \text{element } i \text{ is assigned to bucket } j \\ 0 & otherwise \end{cases}$$

Only true if
$b_{i,j}$ and $b_{i',j}$ are
uncorrelated for any i ≠ i'

$$\mathbb{E}\left[\sum_{i=0}^{n} b_{i,j}\right] = \frac{n}{N}$$

The **expected**
number of elements
in any bucket j

(h(i) can't be related to h(i'))

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Picking a Lookup Function

$$n = \text{number of elements in any bucket}$$

$$N = \text{number of buckets}$$

$$b_{i,j} = \begin{cases} 1 & \textbf{if } \text{element } i \text{ is assigned to bucket } j \\ 0 & otherwise \end{cases}$$

**Expected** Runtime of insert, apply, remove(): $O\left(\dfrac{n}{N}\right)$

**Worst-Case** Runtime of insert, apply, remove(): $O\left(n\right)$

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Hash Functions

- Examples
  - SHA256 ← used by GIT
  - MD5, BCRYPT ← used by unix login, apt
  - MurmurHash3 ← used by Scala
- hash(x) is pseudorandom
  1) hash(x) ~ uniform random value in [0, INT_MAX)
  2) hash(x) always returns the same value
  3) hash(x) uncorrelated with hash(y) for x ≠ y

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Using Hash Functions

- hash(x: Int): Int

  - What about strings?

Arbitrary starting constant
( hash("") )

```
def hashString(str: String): Int = {
  var accumulator: Int = SEED
  for(character <- str){
    accumulator = hash(accumulator * character.toInt)
  }
  return accumulator
}
```

call hash() str.length times

**(simplified, don't actually do exactly this)**

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Hash Functions

- hash(x: Object): Int
    - In Java/Scala, call x.hashCode

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY

# Hash Functions + Buckets

Everything is: $O\left(\dfrac{n}{N}\right)$      Let's call $\alpha = \dfrac{n}{N}$ the load factor.

**Idea:** Make α a constant

Fix an $\alpha_{max}$ and start requiring that $\alpha \leq \alpha_{max}$

**What happens when this constraint is violated?**

©Oliver Kennedy, Eric Mikida, Andrew Hughes
The University at Buffalo, SUNY