# Introduction to Hash Table Data Structure
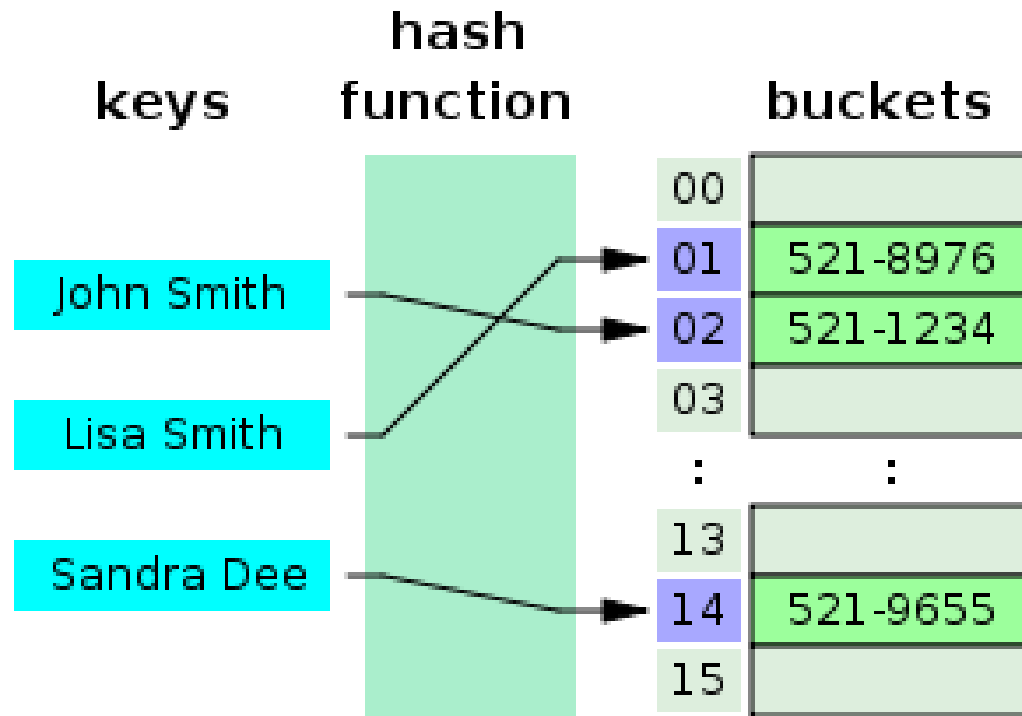
261217 Data Structures for Computer Engineers

Patiwet Wuttisarnwattana, Ph.D.

*patiwet@eng.cmu.ac.th*
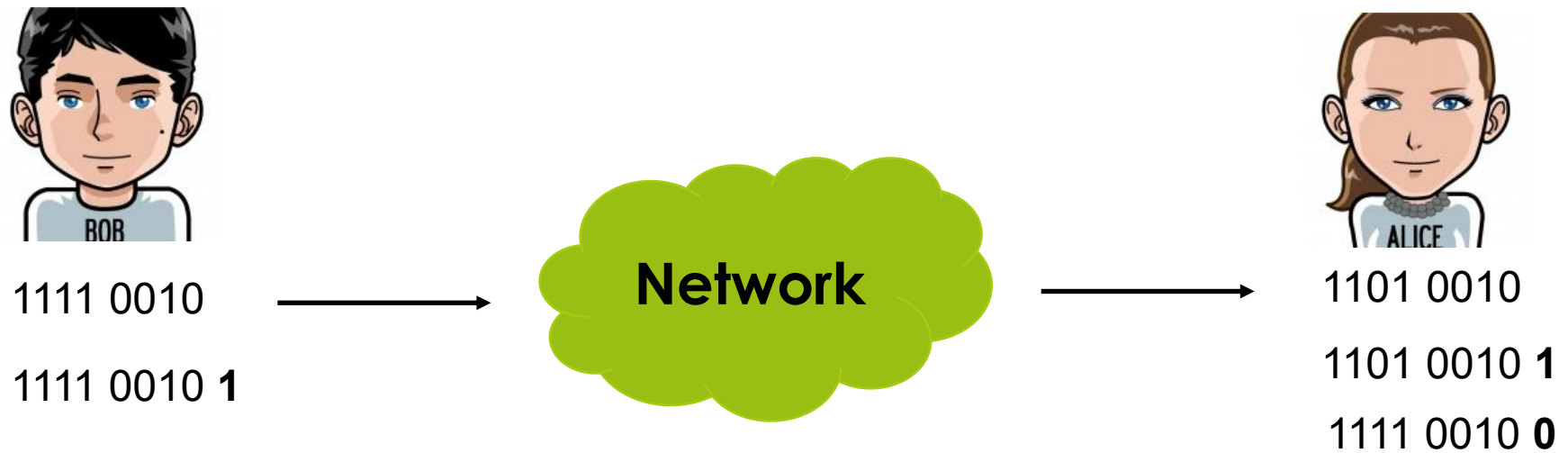
Computer Engineering, Chiang Mai University

# Hash Table Concept



Key Concept is the Hash Function
Hash function can map anything to an integer
The integer is an index of an array (table)

# Hash function Application: Error checking



1111 0010

1111 0010 **1**

**Network**

1101 0010

1101 0010 **1**

1111 0010 **0**

# Bonus Time

$$x = \sum_{i=1}^{12} (14 - i)N_i \quad (\text{mod } 11)$$

$$x = (13N_1 + 12N_2 + 11N_3 + 10N_4 + 9N_5 + 8N_6 + 7N_7 + 6N_8 + 5N_9 + 4N_{10} + 3N_{11} + 2N_{12}) \quad (\text{mod } 11)$$

$$N_{13} = \begin{cases} 1 - x, & \text{if } x \leq 1 \\ 11 - x, & \text{if } x > 1 \end{cases}$$

# Bonus Time

$$x = \sum_{i=1}^{12} (14 - i)N_i \quad (\text{mod } 11)$$

$$x = (13N_1 + 12N_2 + 11N_3 + 10N_4 + 9N_5 + 8N_6 + 7N_7 + 6N_8 + 5N_9 + 4N_{10} + 3N_{11} + 2N_{12}) \quad (\text{mod } 11)$$

$$N_{13} = \begin{cases} 1 - x, & \text{if } x \leq 1 \\ 11 - x, & \text{if } x > 1 \end{cases}$$



1 2345 67890 12 3

แสดงตัวเลขหลักที่ 13
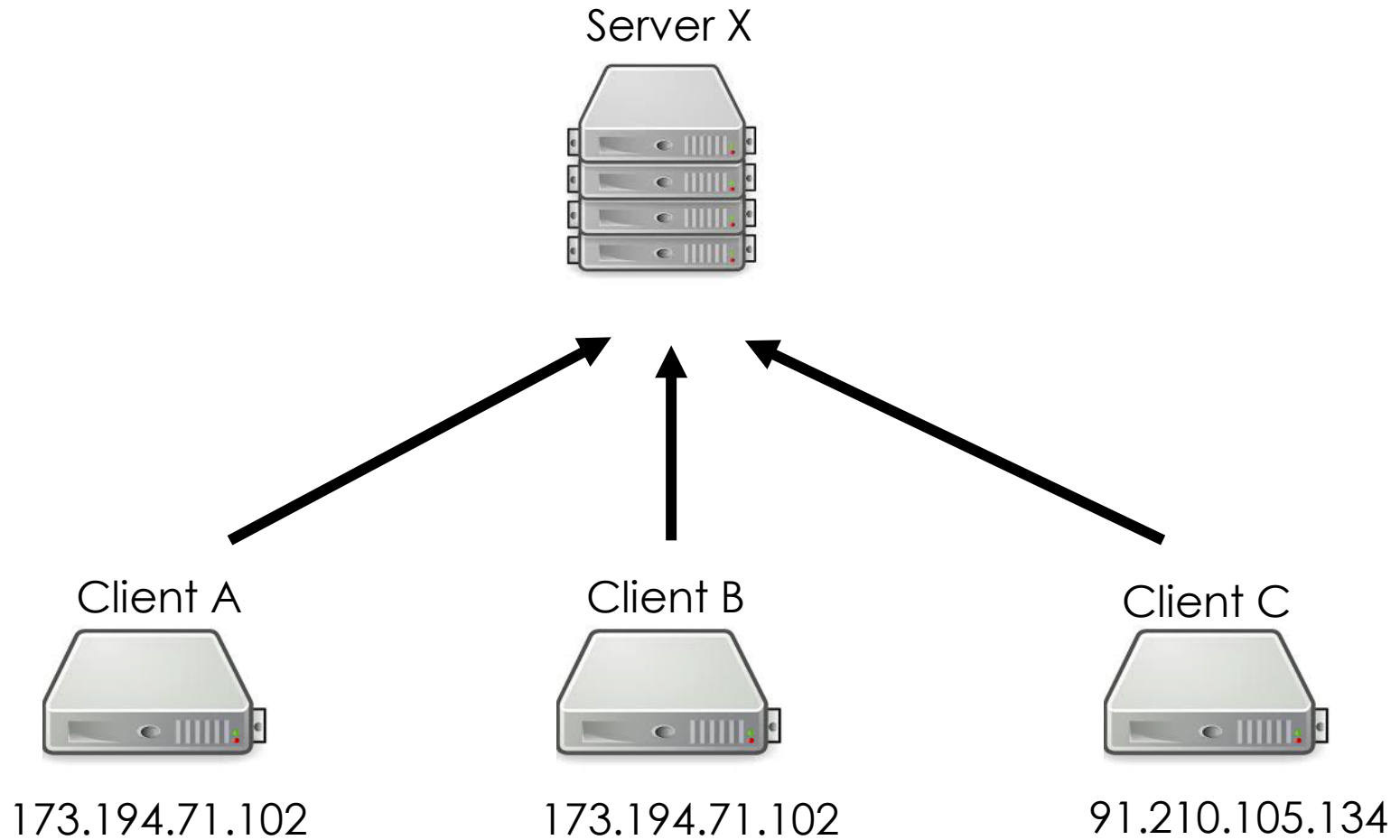
# Hash function Application: Password Storage

Upload the file if the file is not uploaded yet
Guarantee No Duplication

# Web Service

Server X

Client A

173.194.71.102

Client B

173.194.71.102

Client C

91.210.105.134

# Web Service

Server X

**Possible number of IP Addresses**

$2^{32}$ = 4,294,967,296

**Possible number of IPv6 Addresses**

$2^{128}$ = numbers with 39 digits

Client A

Client B

Client C

173.194.71.102

173.194.71.102

91.210.105.134

# Denial of Service Attack

# Access Log

| Date | Time | IP address |
|---|---|---|
| 09 Dec 2015 | 00:45:13 | 173.194.71.102 |
| 09 Dec 2015 | 00:45:15 | 69.171.230.68 |
| ... | ... | ... |
| ... | ... | ... |
| 09 Dec 2015 | 01:45:13 | 91.210.105.134 |

# IP Access List

## IP Access List

Analyse the access log and quickly answer queries: did anybody access the service from this *IP* during the last hour? How many times? How many *IP*s were used to access the service during the last hour?

# Log Processing

- 1h of logs can contain millions of lines

- Too slow to process that for each query

- Keep count: how many times each IP appears in the last 1h of the access log

- **C** is some data structure to store the mapping from IP to counters

- We will learn later how to implement **C**

# Log Processing

**Decrease the counter**

| Time | IP address |
|------|------------|
| 00:45:13 | 173.194.71.102 |
| 00:45:13 | 69.171.230.68 |
| … | … |
| 01:45:13 | 173.194.71.102 |
| 01:45:13 | 91.210.105.134 |

**1 hour ago**

**Now (Need to update)**

**Increase the counter**

# Main Loop

$log$ - array of log lines $(time, IP)$
$C$ - mapping from IPs to counters
$i$ - first unprocessed log line
$j$ - first line in current 1h window
$i \leftarrow 0$
$j \leftarrow 0$
$C \leftarrow \emptyset$
Each second
    UpdateAccessList$(log, i, j, C)$

## UpdateAccessList$(log, i, j, C)$

```
while log[i].time ≤ Now():
   C[log[i].IP] ← C[log[i].IP] + 1
   i ← i + 1
while log[j].time ≤ Now() − 3600:
   C[log[j].IP] ← C[log[j].IP] − 1
   j ← j + 1
```

$$\text{while } log[i].time \leq Now():$$
$$\quad C[log[i].IP] \leftarrow C[log[i].IP] + 1$$
$$\quad i \leftarrow i + 1$$
$$\text{while } log[j].time \leq Now() - 3600:$$
$$\quad C[log[j].IP] \leftarrow C[log[j].IP] - 1$$
$$\quad j \leftarrow j + 1$$

## AccessedLastHour$(IP, C)$
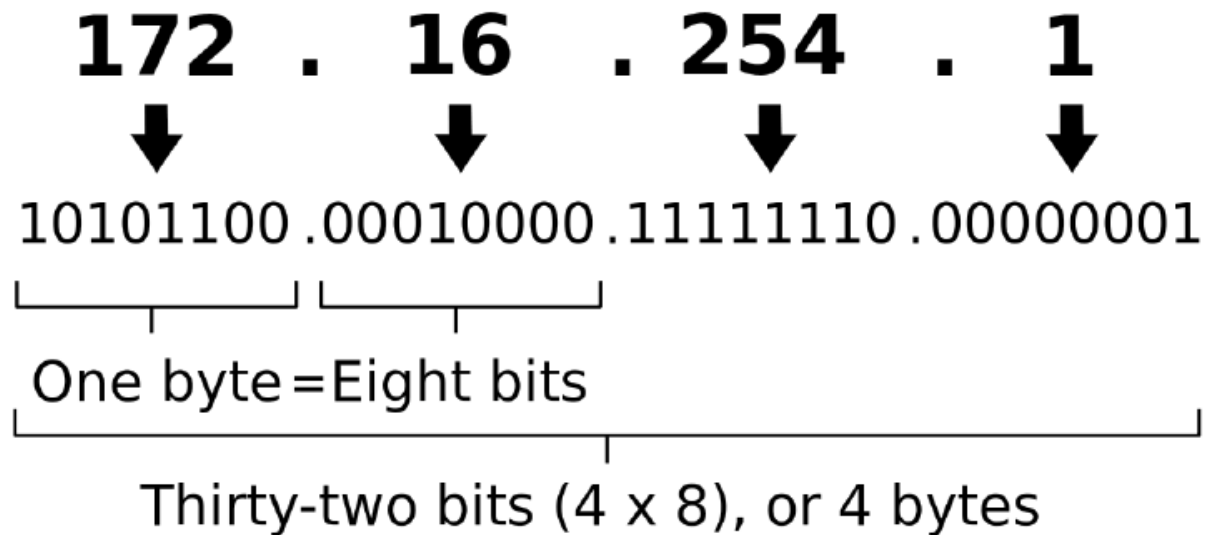
$$\text{return } C[IP] > 0$$

The Question is How to implement the mapping C?

# Direct Addressing

- Need a data structure for **C**

- There are $2^{32}$ different IP(v4) addresses
  - $2^{32}$ = 4,294,967,296

- Convert IP to 32-bit integer

- Create an integer array `A` of size $2^{32}$

- Use `A[int(IP)]` as `C[IP]`

# int(IP)

An IPv4 address (dotted-decimal notation)

**172** . **16** . **254** . **1**

10101100 . 00010000 . 11111110 . 00000001

One byte = Eight bits

Thirty-two bits (4 x 8), or 4 bytes

- int(0.0.0.1) = 1
- int(172.16.254.1) = 2,886,794,753
- int(69.171.230.68) = 1,168,893,508

## int($IP$)

return $IP[1]\cdot 2^{24} + IP[2]\cdot 2^{16} + IP[3]\cdot 2^{8} + IP[4]$

## UpdateAccessList($log, i, j, A$)

```
while log[i].time ≤ Now():
    A[int(log[i].IP)] ← A[int(log[i].IP)] + 1
    i ← i + 1
while log[j].time ≤ Now() − 3600:
    A[int(log[j].IP)] ← A[int(log[j].IP)] − 1
    j ← j + 1
```

Is the server accessed by the IP in the last hour?

$$AccessedLastHour(IP)$$

$$\text{return } A[\text{int}(IP)] > 0$$

# Big O Analysis of Direct Addressing Implementation

- UpdateAccessList is O(1) per log line

- AccessedLastHour is O(1)

- But need $2^{32}$ memory (4GB) even for few IPs

- IPv6: $2^{128}$ will not fit in memory

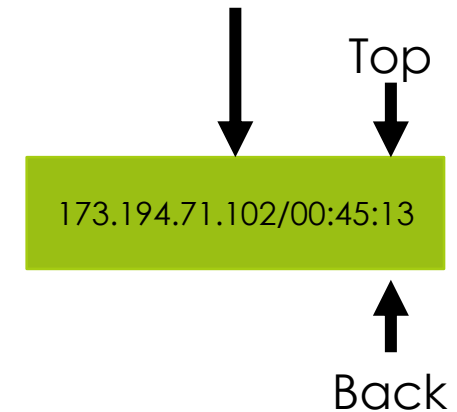- In general: O(**U**) memory, **U** = number of possible IPs

# List-based Mapping

- Direct addressing requires too much memory

- Let's store only active IPs

- Store them in a list

- Store only last occurrence of each IP

- Keep the order of occurrence

# Access Log

Current Time: 00:45:15

**List**

Top

173.194.71.102/00:45:13

Back

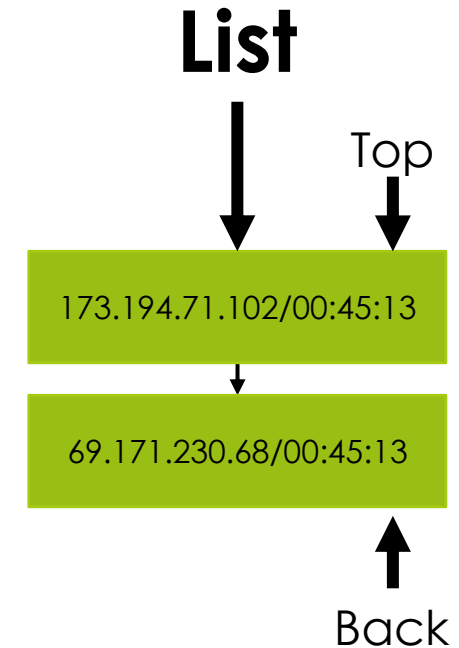| Time | IP address |
|---|---|
| 00:45:13 | 173.194.71.102 |
| 00:45:13 | 69.171.230.68 |
| 01:00:00 | 69.171.230.68 |
| 01:45:13 | 173.194.71.102 |
| 01:45:13 | 91.210.105.134 |

AccessedLastHour(69.171.230.68)?          AccessCountLastHour(69.171.230.68)?

# Access Log

Current Time: 00:45:16

| Time | IP address |
|------|------------|
| 00:45:13 | 173.194.71.102 |
| 00:45:13 | 69.171.230.68 |
| 01:00:00 | 69.171.230.68 |
| 01:45:13 | 173.194.71.102 |
| 01:45:13 | 91.210.105.134 |

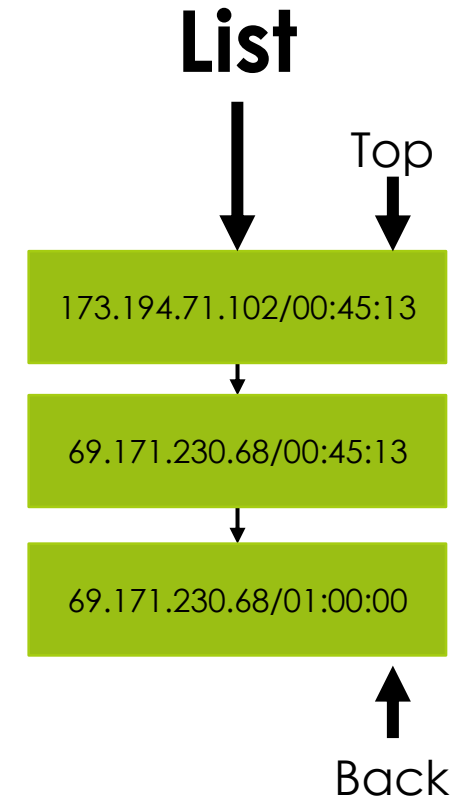**List**

Top

173.194.71.102/00:45:13

69.171.230.68/00:45:13

Back

AccessedLastHour(69.171.230.68)?          AccessCountLastHour(69.171.230.68)?

# Access Log

Current Time: 01:00:05

| Time | IP address |
|------|------------|
| 00:45:13 | 173.194.71.102 |
| 00:45:13 | 69.171.230.68 |
| 01:00:00 | 69.171.230.68 |
| 01:45:13 | 173.194.71.102 |
| 01:45:13 | 91.210.105.134 |

**List**

Top

173.194.71.102/00:45:13

69.171.230.68/00:45:13

69.171.230.68/01:00:00

Back

AccessedLastHour(69.171.230.68)?          AccessCountLastHour(69.171.230.68)?

# Access Log

Current Time: 01:45:15

**List**

| Time | IP address |
|------|-----------|
| 00:45:13 | 173.194.71.102 |
| 00:45:13 | 69.171.230.68 |
| 01:00:00 | 69.171.230.68 |
| 01:45:13 | 173.194.71.102 |
| 01:45:13 | 91.210.105.134 |

Top

69.171.230.68/00:45:13

69.171.230.68/01:00:00

Back

AccessedLastHour(69.171.230.68)?    AccessCountLastHour(69.171.230.68)?

# Access Log

Current Time: 01:45:16

**List**

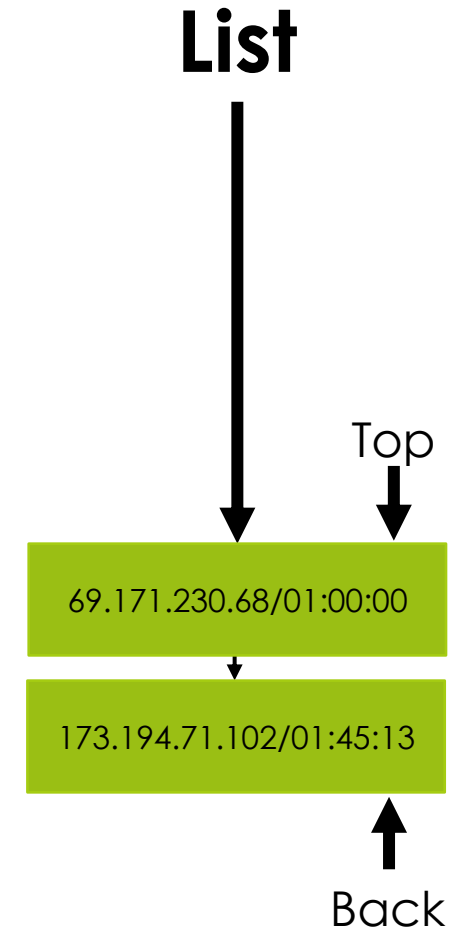| Time | IP address |
|------|------------|
| 00:45:13 | 173.194.71.102 |
| 00:45:13 | 69.171.230.68 |
| 01:00:00 | 69.171.230.68 |
| 01:45:13 | 173.194.71.102 |
| 01:45:13 | 91.210.105.134 |

Top

69.171.230.68/01:00:00

Back

AccessedLastHour(69.171.230.68)?          AccessCountLastHour(69.171.230.68)?

# Access Log

Current Time: 01:45:17

| Time | IP address |
|------|------------|
| 00:45:13 | 173.194.71.102 |
| 00:45:13 | 69.171.230.68 |
| 01:00:00 | 69.171.230.68 |
| 01:45:13 | 173.194.71.102 |
| 01:45:13 | 91.210.105.134 |

**List**

Top

69.171.230.68/01:00:00

173.194.71.102/01:45:13

Back

# Access Log

Current Time: 01:45:18

| Time | IP address |
|---|---|
| 00:45:13 | 173.194.71.102 |
| 00:45:13 | 69.171.230.68 |
| 01:00:00 | 69.171.230.68 |
| 01:45:13 | 173.194.71.102 |
| 01:45:13 | 91.210.105.134 |

**List**

Top

69.171.230.68/01:00:00

173.194.71.102/01:45:13

91.210.105.134/01:45:13

Back

# UpdateAccessList(*log*, *i*, *L*)

**UpdateAccessList(*log*, *i*, *L*)**

```
while log[i].time ≤ Now()
     L.Append(log[i])
     i ← i + 1
while L.Top().time ≤ Now() - 3600
     L.Pop()
```

# AccessedLastHour and AccessCountLastHour

**AccessedLastHour(*IP, L*)**

return *L*.FindByIP(*IP*) ≠ NULL

**AccessCountLastHour(*IP, L*)**

return *L*.CountIP(*IP*)

# Big O Analysis of List-based Implementation

- n is number of active IPs

- Memory usage is O(n)

- `L.Append, K.Top, L.pop` are O(1)

- `UpdateAccessList` is O(1) per log line

- `L.FindByIP` and `L.CountIP` are O(n)

- `AccessedLastHour` and `AccessCountLastHour` are O(n)

# Encoding (Hashing) IPs

- Encode/Hash IPs with small numbers

- For example, numbers from 0 to 999

- Different codes for currently active IPs

# Hash Function

## Definition

For any set of objects $S$ and any integer $m > 0$, a function $h : S \to \{0, 1, \ldots, m - 1\}$ is called a hash function.

## Definition

$m$ is called the cardinality of hash function $h$.

# Bonus Time!!!

# Desirable Properties

- **h** should be fast to compute

- Different values for different objects

- Direct addressing with $O(\mathbf{m})$ memory

- Want small cardinality **m**

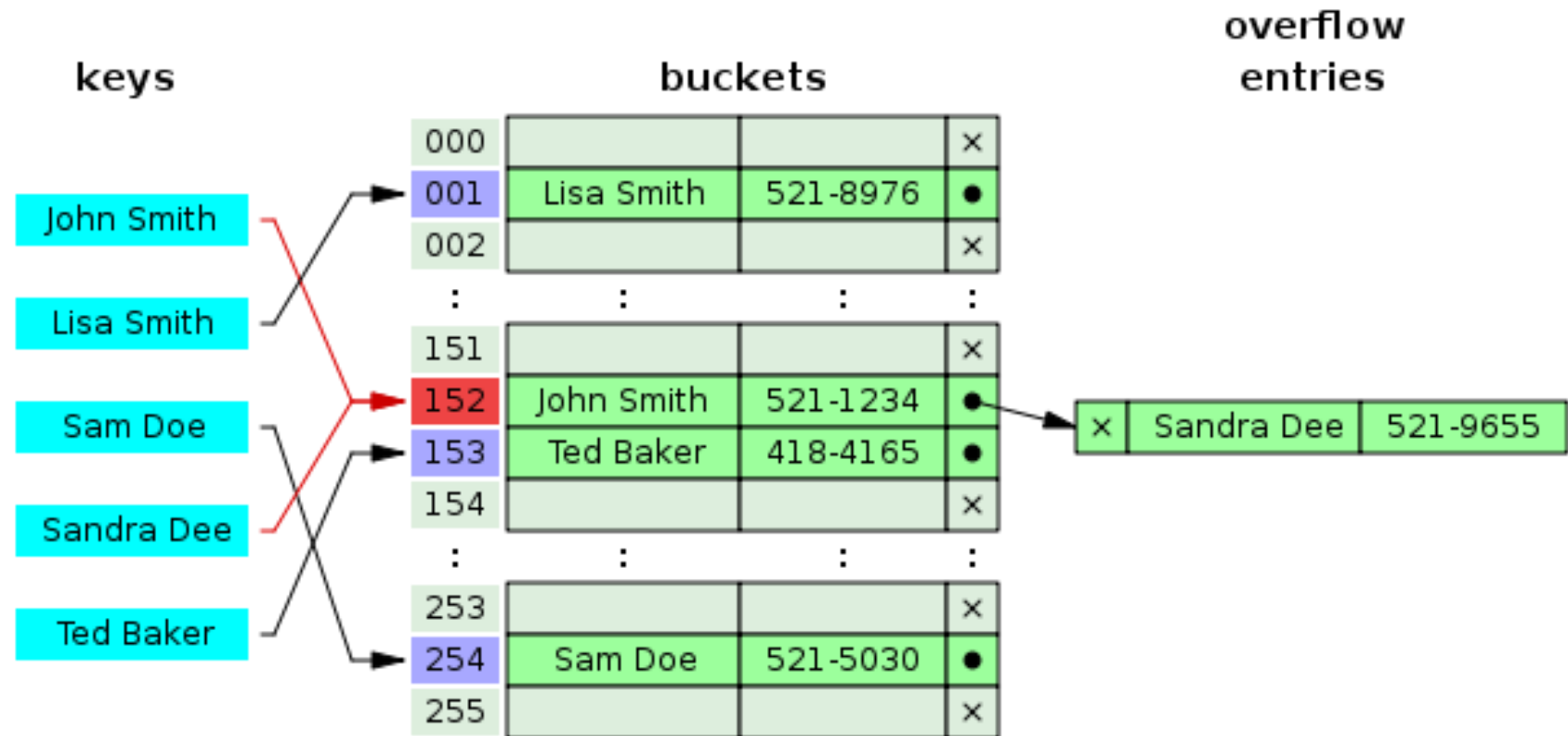- Impossible to have all different values if number of objects in *the universe* is more than **m**
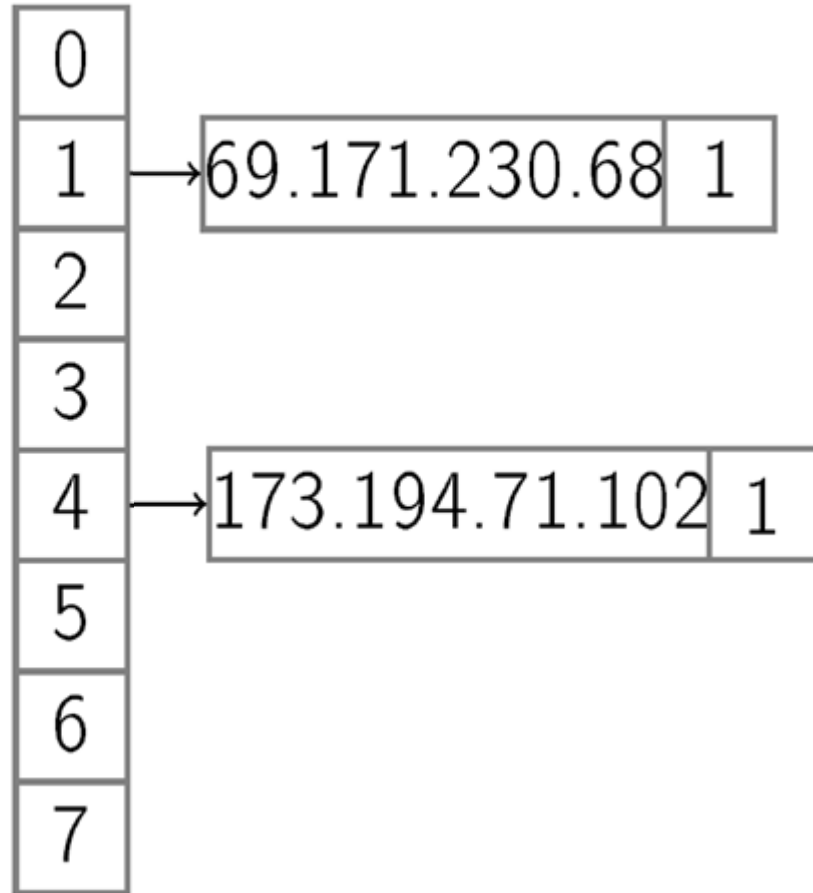
# Collisions

## Definition

When $h(o_1) = h(o_2)$ and $o_1 \neq o_2$, this is a collision.

**Separate Chaining**

# Back to our example

| | |
|---|---|
| 0 | |
| 1 | → 69.171.230.68 \| 1 |
| 2 | |
| 3 | |
| 4 | → 173.194.71.102 \| 1 |
| 5 | |
| 6 | |
| 7 | |

$h(173.194.71.102) = 4$

$h(69.171.230.68) = 1$

$h(173.194.71.102) = 4$

# There is no collision

| | |
|---|---|
| 0 | |
| 1 | → 69.171.230.68 \| 1 |
| 2 | |
| 3 | |
| 4 | → 173.194.71.102 \| 2 → 91.210.105.134 \| 1 |
| 5 | |
| 6 | |
| 7 | |

$$h(173.194.71.102) = 4$$

$$h(69.171.230.68) = 1$$

$$h(173.194.71.102) = 4$$

$$h(91.210.105.134) = 4$$

Map Data Structure

Dictionary Data Structure

Associative Array Data Structure

are Hash Table

# Map (Dictionary, Associative Array)

- Arrays store items as an ordered collection and you can access thems with an index number.

- Maps store items in (Key, Value) pairs that you can access values by the corresponding keys (usually Strings)
    - Filename → location of the file on disk
    - Student ID → student name
    - Contact name → contact phone number
  - The first object is called "Key" or O
  - The second object is called "Value" or v
  - Map can have the following operations `HasKey(O)`, `Get(O), Set(O,v)`

# Map Data Structure

Map217

```
{
    {"580610615", "Kanokwan Pinthong"},
    {"580610616", "Kawewut Chujit"},
    {"580610618", "Kittitorn Rakpanyakeaw"}
}
```

- HasKey("580610615")
- HasKey("580610617")
- Get("580610618")
- Set("580610618", "Kittitat Boonkarn")
- Get("580610618")

# Demo: Java HashMap

https://www.w3schools.com/java/java_hashmap.asp

Map Data Structure

Dictionary Data Structure

Associative Array Data Structure

can be implemented using Hash Table

$$h : S \to \{0, 1, \ldots, m - 1\}$$

$$O, O' \in S$$

$$v, v' \in V$$

$A \leftarrow$ array of $m$ lists (chains) of pairs $(O, v)$

HasKey($O$)

$L \leftarrow A[h(O)]$

```
for (O', v') in L:
    if O' == O:
        return true
return false
```

## Get($O$)

$$L \leftarrow A[h(O)]$$

```
for (O', v') in L:
    if O' == O:
        return v'
return n/a
```

## Set($O, v$)

```
L ← A[h(O)]
for p in L:
    if p.O == O:
        p.v ← v
        return
L.Append(O, v)
```

# Runtime Analysis

## Lemma

Let $c$ be the length of the longest chain in $A$. Then the running time of HasKey, Get, Set is $O(c + 1)$.

## Lemma

Let $n$ be the number of different keys $O$ currently in the map and $m$ be the cardinality of the hash function. Then the memory consumption for chaining is $\Theta(n + m)$.

# Set

## Definition

Set is a data structure with methods
$\mathrm{Add}(O)$, $\mathrm{Remove}(O)$, $\mathrm{Find}(O)$.

## Examples

- IPs accessed during last hour
- Students on campus
- Keywords in a programming language

# Implementing Set

$$h : S \rightarrow \{0, 1, \ldots, m - 1\}$$
$$O, O' \in S$$

$A \leftarrow$ array of $m$ lists (chains) of objects $O$

Find($O$)

```
L ← A[h(O)]
for O' in L:
    if O' == O:
        return true
return false
```

# Implementing Set

$\text{Add}(O)$

$L \leftarrow A[h(O)]$
```
for O' in L:
    if O' == O:
        return
```
$L.\text{Append}(O)$

# Implementing Set

Remove($O$)

if not Find($O$):
   return
$L \leftarrow A[h(O)]$
$L$.Erase($O$)

# Hash Table

**Definition**

An implementation of a set or a map using hashing is called a hash table.

# Programming Languages

Set:

- `unordered_set` in C++

- `HashSet` in Java

- `set` in Python

Map:

- `unordered_map` in C++

- `HashMap` in Java

- `dict` in Python

# Conclusion

- Chaining is a technique to implement a hash table

- Memory consumption is $O(n + m)$

- Operations works in time $O(c + 1)$

- How to make both m and c small?