

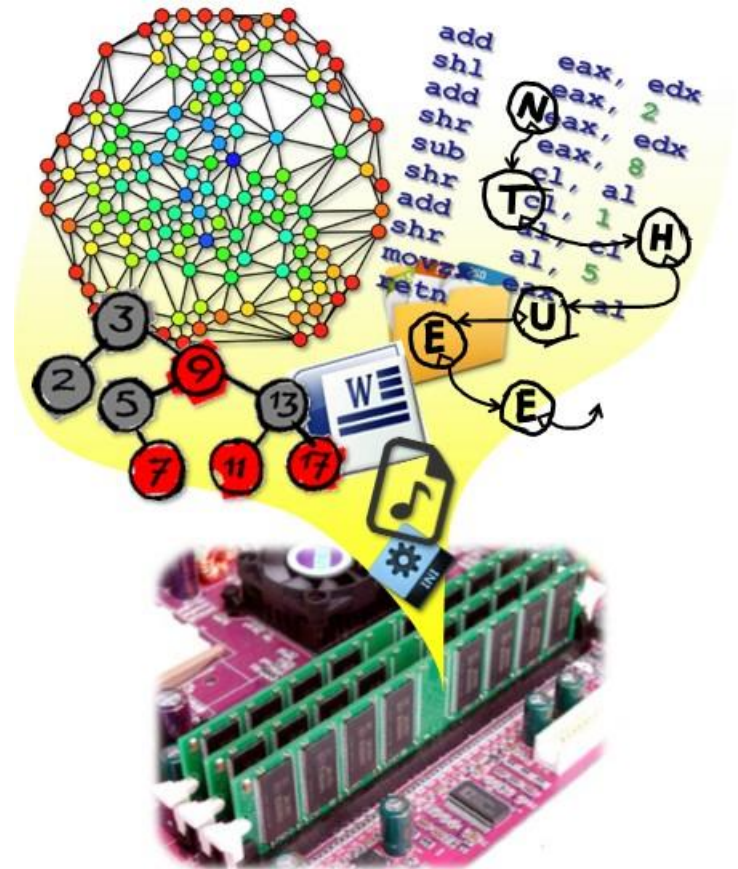
Data Structures

CH8 Hashing

Prof. Ren-Shuo Liu

NTHU EE

Spring 2019

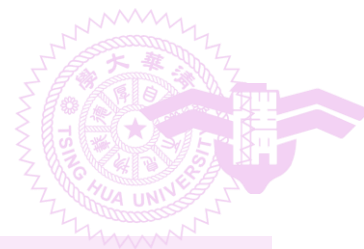




Outline

- 8.1 Introduction
- (8.2 Static hashing)
- (8.3 Dynamic hashing)
- 8.4 Bloom filters
- Security hash

Registration Division Example



請大家向註冊組承辦人
查詢學期成績



| 承辦人 | 分機 / Email |
|-----|----------------------|
| 陳OO | 31300 / chen@nthu... |
| 郭OO | 31301 / kuo@nthu... |
| 李OO | 31302 / li@nthu... |
| 林OO | 31303 / lin@nthu.. |
| 王OO | 31304 / wang@nthu... |

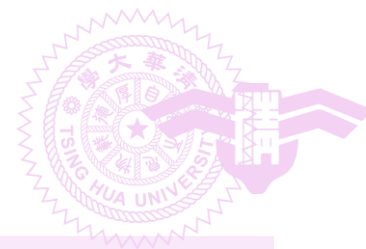
Registration Division Example



全校多數都打電話、
寄email給第一位



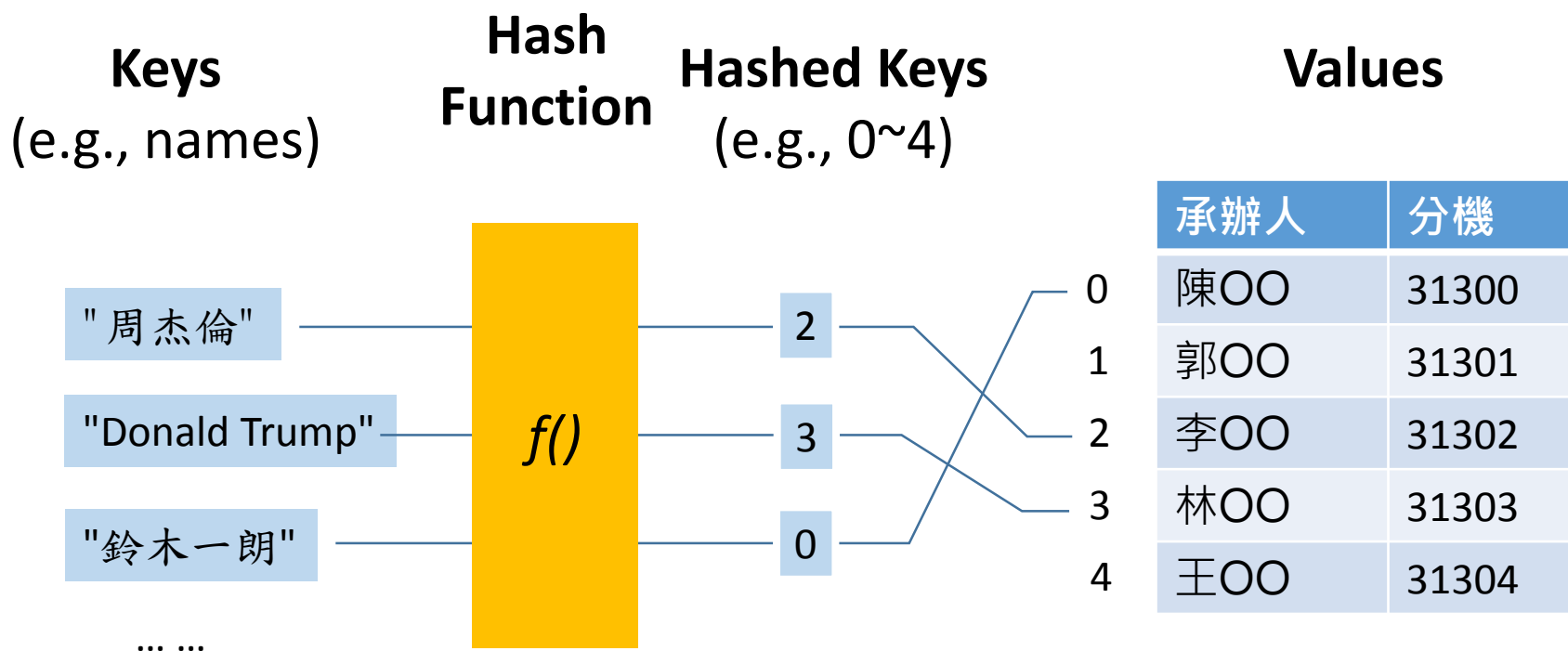
| 承辦人 | 分機 / Email |
|-----|----------------------|
| 陳OO | 31300 / chen@nthu... |
| 郭OO | 31301 / kuo@nthu... |
| 李OO | 31302 / li@nthu... |
| 林OO | 31303 / lin@nthu.. |
| 王OO | 31304 / wang@nthu... |

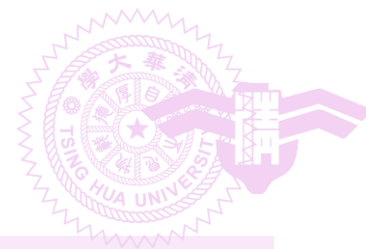


Hash Concepts

- **Hash function**

- Any deterministic function that can map data of arbitrary size (original keys) to data of a desired fixed size (hashed keys)

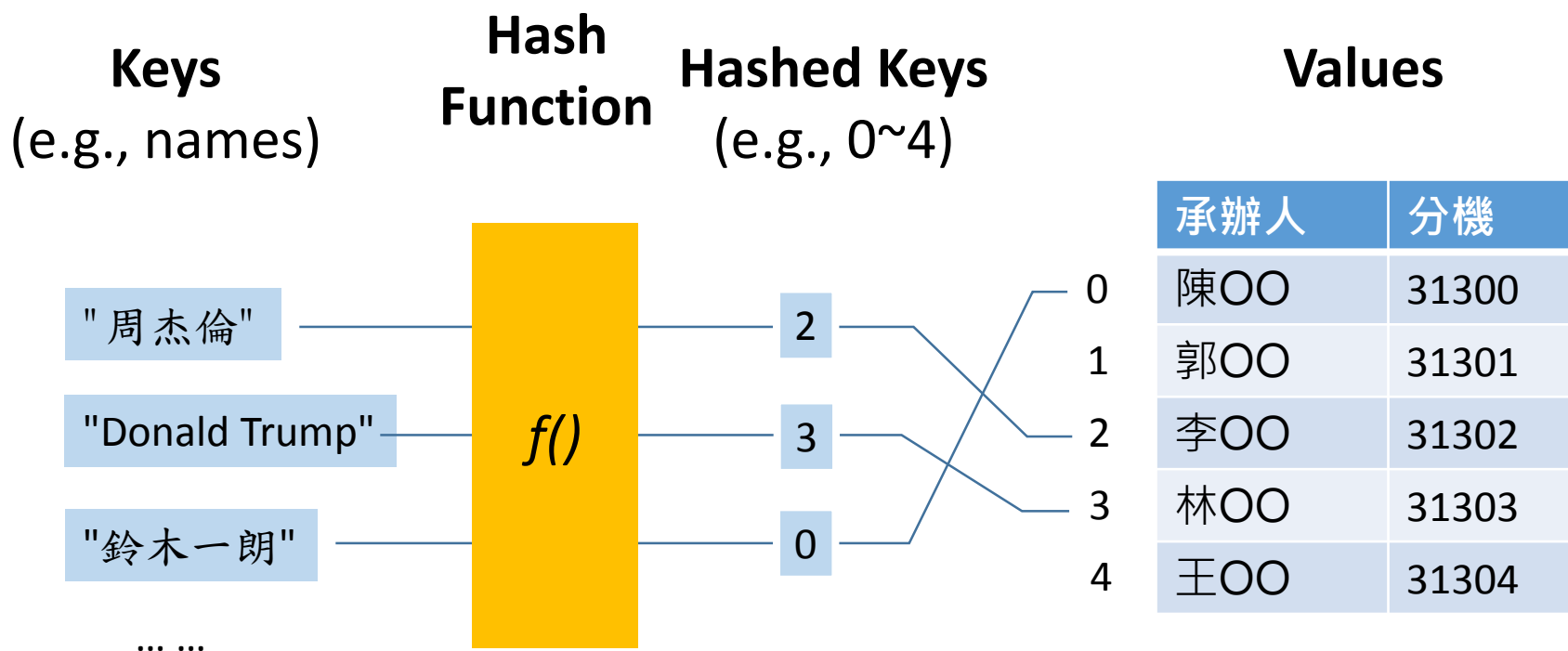




Hash Concepts

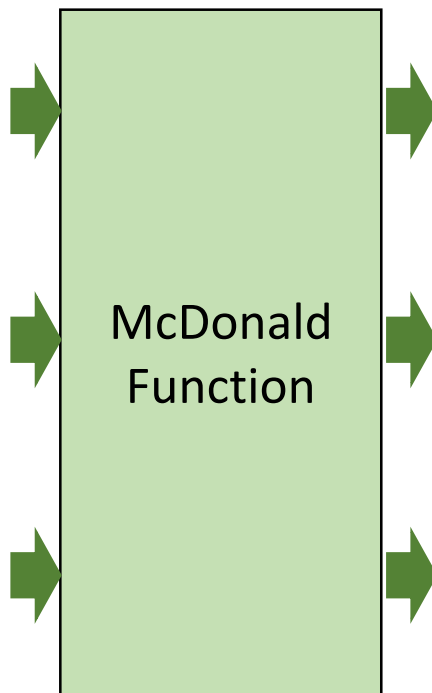
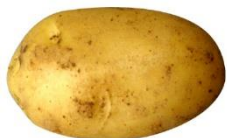
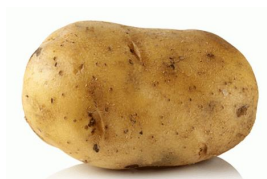
- **Hash function**

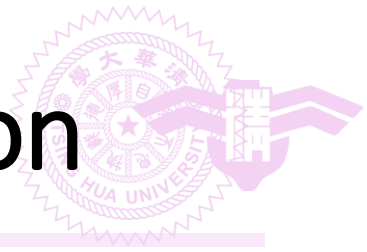
- It shuffles the order of mapping
- But it is deterministic



Hash in Cooking

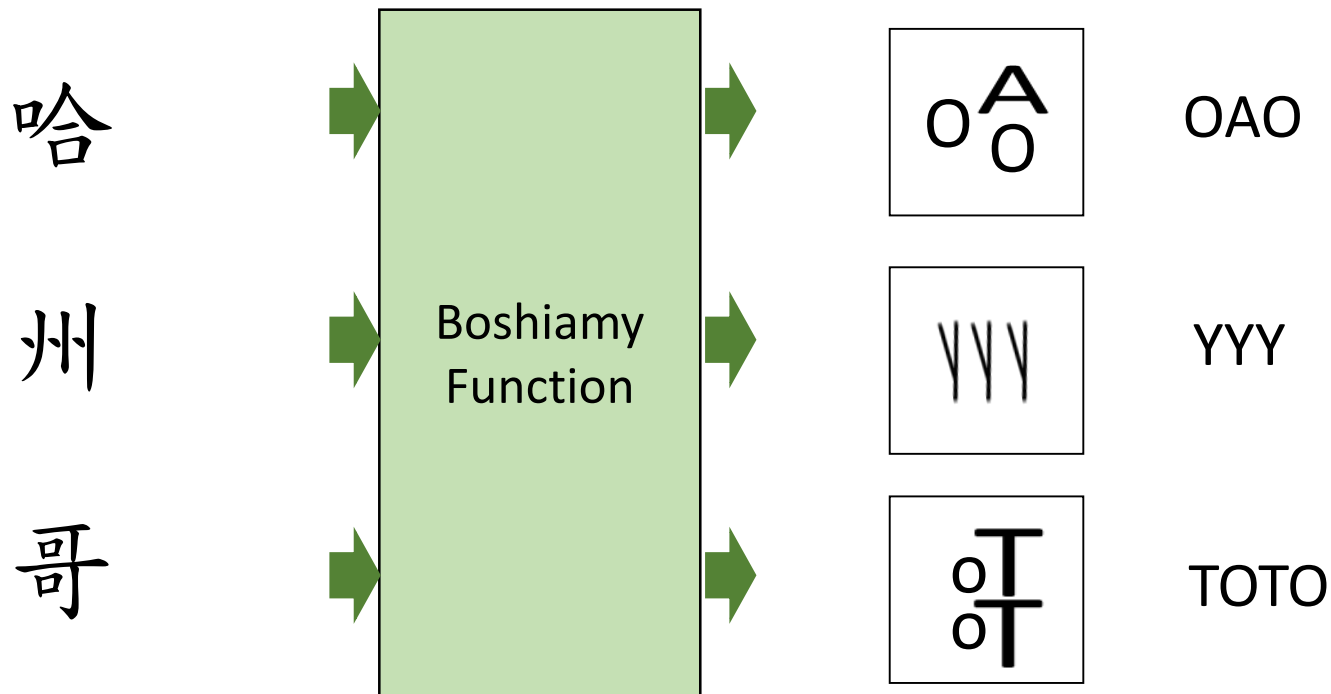
- Hash: "chop and mix foods"
- Example: hash browns (薯餅)





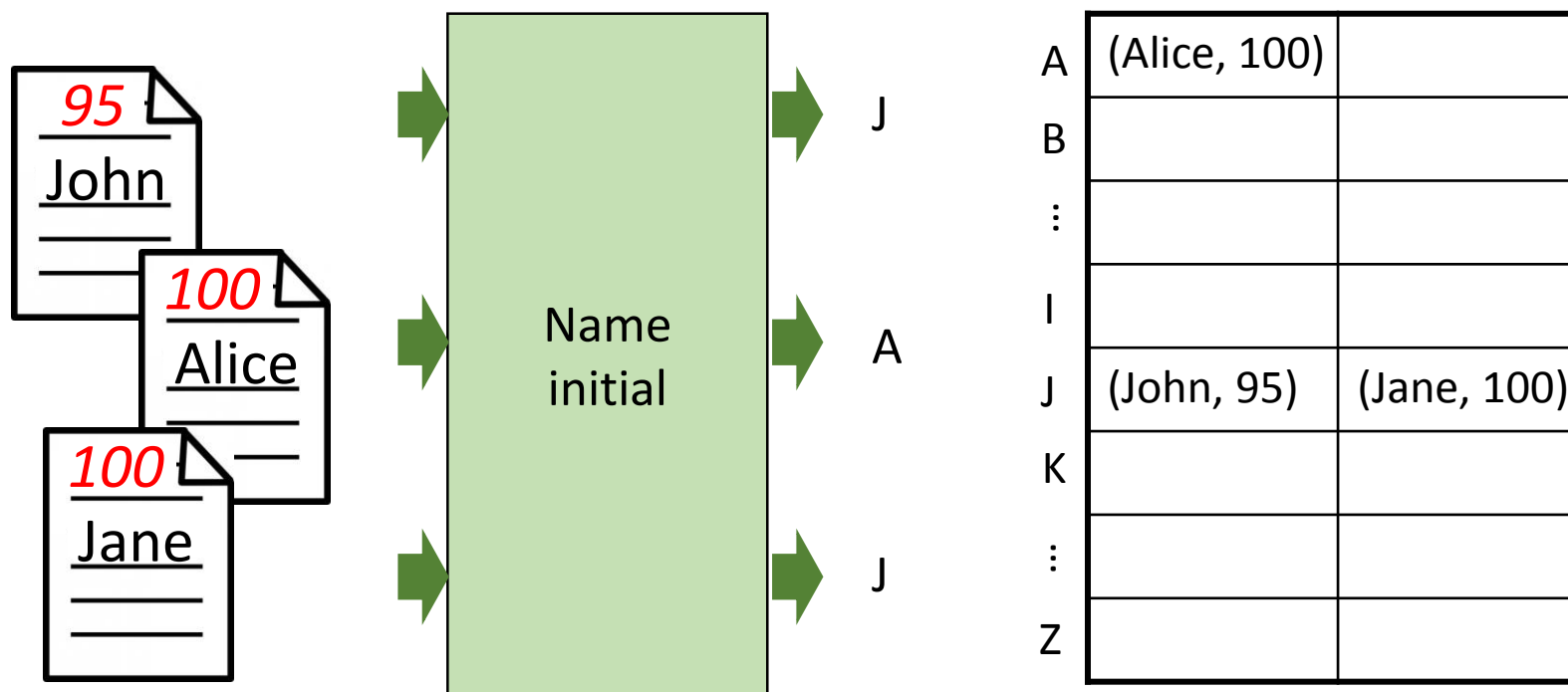
Hash in Chinese Decomposition

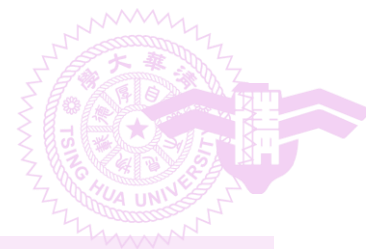
- Decompose Chinese characters into keyboard strokes
 - Facilitate Chinese input
- Example: the Boshiamy (嘸蝦米) decomposition scheme



Hash in a Storing Data

- Example: Storing students' grades according to their name initial letters

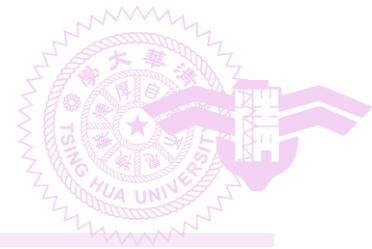




Advantages of Hashing

- Inserting, deleting, and searching can be $\sim O(1)$ time
 - Hash function computation is designed in $O(1)$
 - Indexing the corresponding bucket in the table is $O(1)$
 - Searching all slots in a bucket for a key is also $O(1)$
 - The number of slots is independent of the number of pairs stored in the table

| | | |
|---|--------------|-------------|
| A | (Alice, 100) | |
| B | (Bob, 80) | (Ben, 70) |
| ⋮ | | |
| I | (Irene, 85) | |
| J | (John, 95) | (Jane, 100) |
| K | (Ken, 75) | |
| ⋮ | | |
| Z | (Zoe, 80) | |

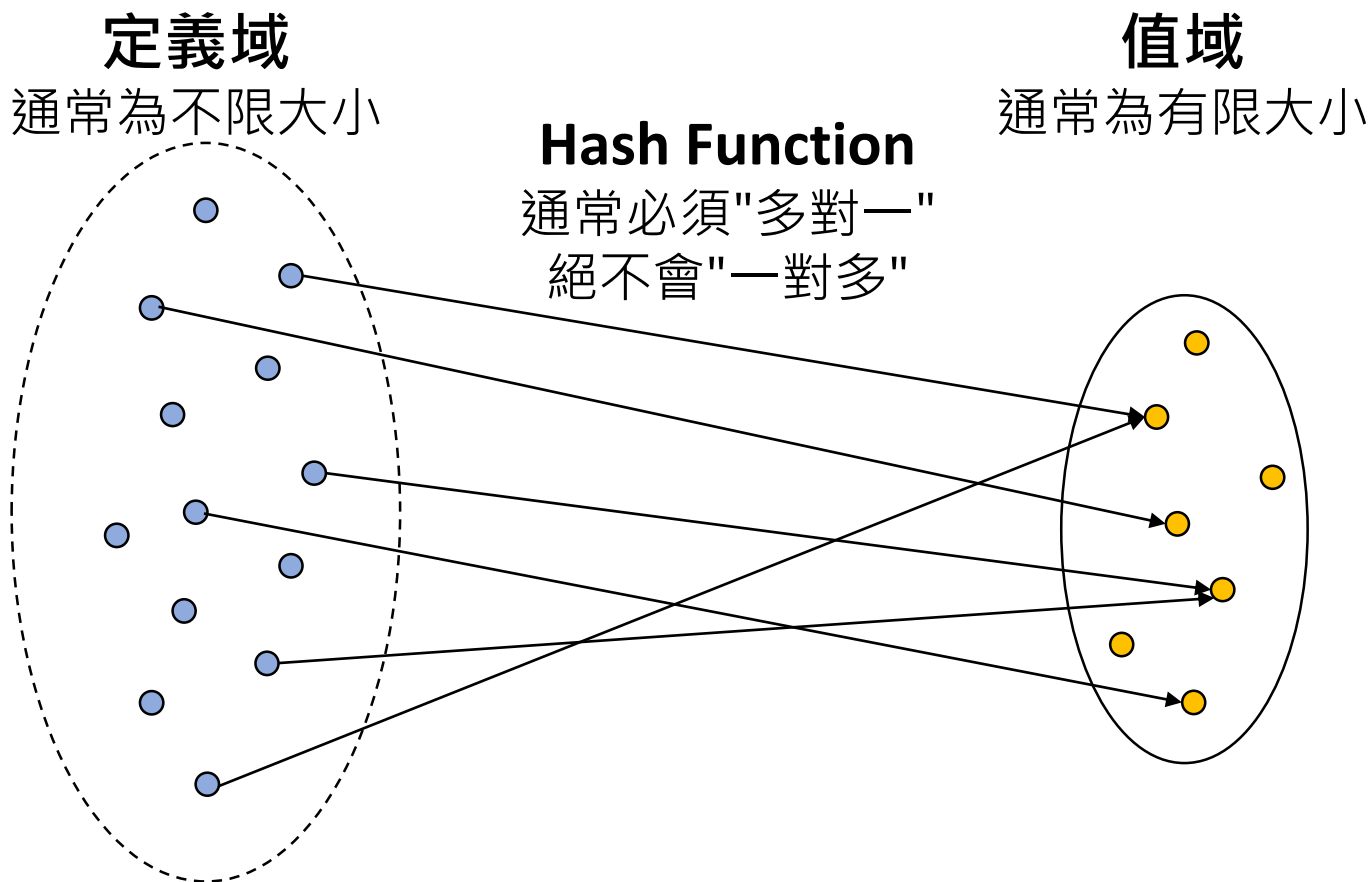


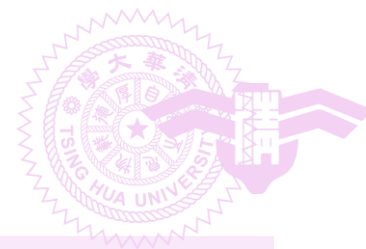
Hashing

- A pair with a key k is stored in a hash table ht
- Key parameters
 - b buckets in ht
 - $h(k)$ is the home bucket of a key k
 - s slots per bucket
 - T possible different keys
 - n stored pairs in ht

| Slots | |
|---------|--------------------------|
| Buckets | A (Alice, 100) |
| | B (Bob, 80) (Ben, 70) |
| | ⋮ |
| | I (Irene, 85) |
| | J (John, 95) (Jane, 100) |
| | K (Ken, 75) |
| | ⋮ |
| | Z (Zoe, 80) |

Hash Function

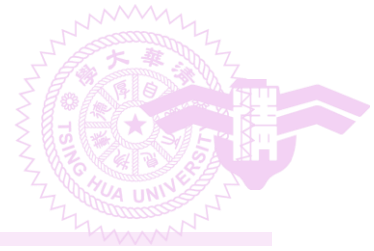




Hash Function

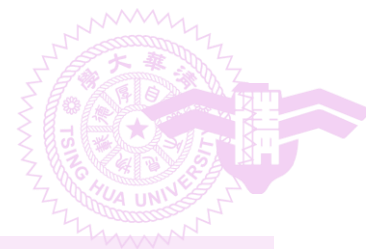
- Good hash functions reduce the chance of collisions
- Enlarging hash table size can also reduce collisions
 - At the cost of memory size
- Ideal hash functions
 - Uniform hash function
 - Easy to compute

| Slots | |
|---------|--------------------------|
| Buckets | A (Alice, 100) |
| | B (Bob, 80) (Ben, 70) |
| | ⋮ |
| | I (Irene, 85) |
| | J (John, 95) (Jane, 100) |
| | K (Ken, 75) |
| | ⋮ |
| | Z (Zoe, 80) |



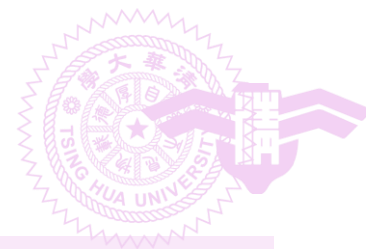
Hash Functions

- Classical examples
 - Modulo (division)
 - Mid-square
 - Folding
 - String-to-integer conversion
- We can design our own hash functions



Modulo (Division)

- Most widely used hash function in practice
- Procedure
 - $h(k) = k \% D$
- Selection of D
 - $D \leq$ the number of buckets
 - D would better be an odd number
 - Even divisor D always maps even keys to even buckets and odd keys to odd buckets
 - Real-world data tend to have a bias toward either odd or even keys
 - It would be even desirable if D can be a prime number or a number having no prime factors smaller than 20

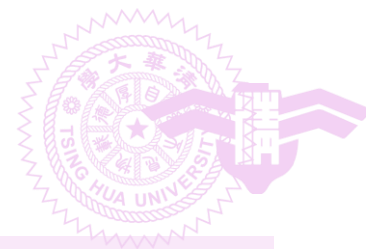


Mid-Square

- $h(k)$ = some middle r bits of the square of k
 - The number of bucket is equal to 2^r
- Example

| k | k^2 | | $h(k)$ |
|---|-------|-------------------|--------|
| 0 | 0 | 00 <u>00</u> 0000 | 0 |
| 1 | 1 | 00 <u>00</u> 0001 | 0 |
| 2 | 4 | 00 <u>00</u> 0100 | 1 |
| 3 | 9 | 00 <u>00</u> 1001 | 2 |
| 4 | 16 | 00 <u>01</u> 0000 | 4 |
| 5 | 25 | 00 <u>01</u> 1001 | 6 |
| 6 | 36 | 00 <u>10</u> 0100 | 9 |
| 7 | 49 | 00 <u>11</u> 0001 | 12 |

| k | k^2 | | $h(k)$ |
|----|-------|-------------------|--------|
| 8 | 64 | 01 <u>00</u> 0000 | 0 |
| 9 | 81 | 01 <u>01</u> 0001 | 4 |
| 10 | 100 | 01 <u>10</u> 0100 | 9 |
| 11 | 121 | 01 <u>11</u> 1001 | 14 |
| 12 | 144 | 10 <u>01</u> 0000 | 4 |
| 13 | 169 | 10 <u>10</u> 1001 | 10 |
| 14 | 196 | 11 <u>00</u> 0100 | 1 |
| 15 | 225 | 11 <u>10</u> 0001 | 8 |



Folding

- Partition the key into several parts and add them together
 - Two strategies: **shift folding** and **folding at the boundary**

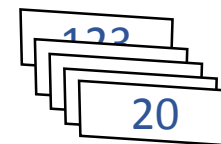
- Example

- $k = 12320324111220 =$

| | | | | |
|-----|-----|-----|-----|----|
| 123 | 203 | 241 | 112 | 20 |
|-----|-----|-----|-----|----|

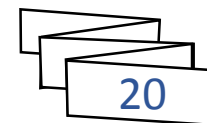
- Shift folding

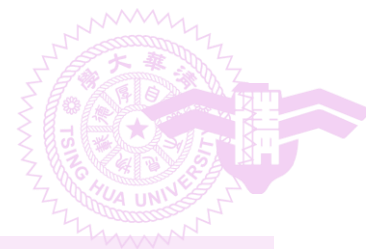
$$h(k) = \sum \begin{array}{|c|c|c|c|c|} \hline 123 & 203 & 241 & 112 & 20 \\ \hline \end{array} = 699$$



- Folding at the boundary

$$h(k) = \sum \begin{array}{|c|c|c|c|c|} \hline 123 & 302 & 241 & 211 & 20 \\ \hline \end{array} = 897$$



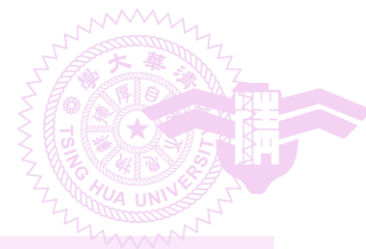


String-to-Integer Conversion

- Useful when keys are strings
- Procedure
 - Treat every n character as an 8n-bit integer
 - ASCII represents a character using 8 bits

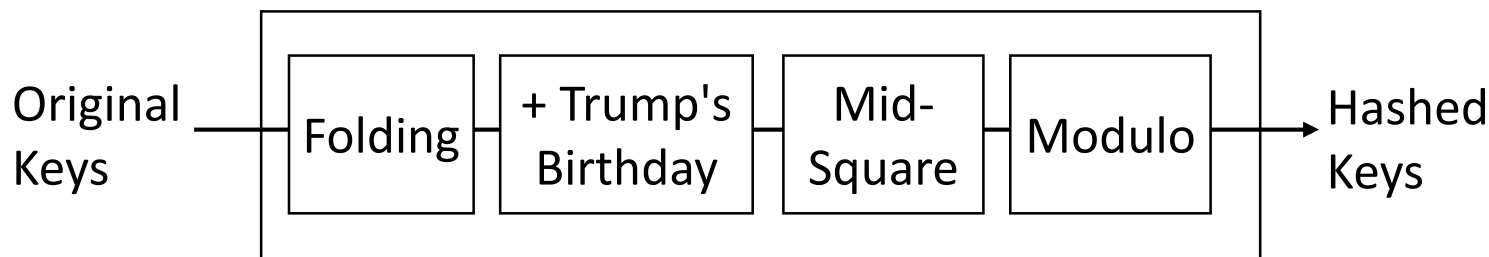
| | | | | |
|----------------|----------|----------|----------|----------|
| Characters: | h | o | p | e |
| ASCII Values: | 104 | 111 | 112 | 101 |
| Binary Values: | 01101000 | 01101111 | 01110000 | 01100101 |

- Add all integers together to obtain the overall value
- Adopt the aforementioned hash functions (modulo, folding...)

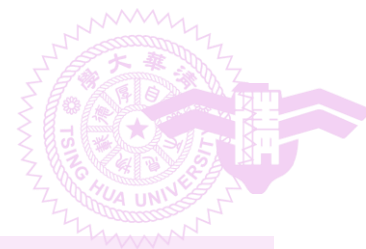


Design Our Own Hash

- Recall that
 - Hash function is **any** deterministic function that can map data of arbitrary size (original keys) to data of a desired fixed size (hashed keys)
- So of course we can design a hash like this

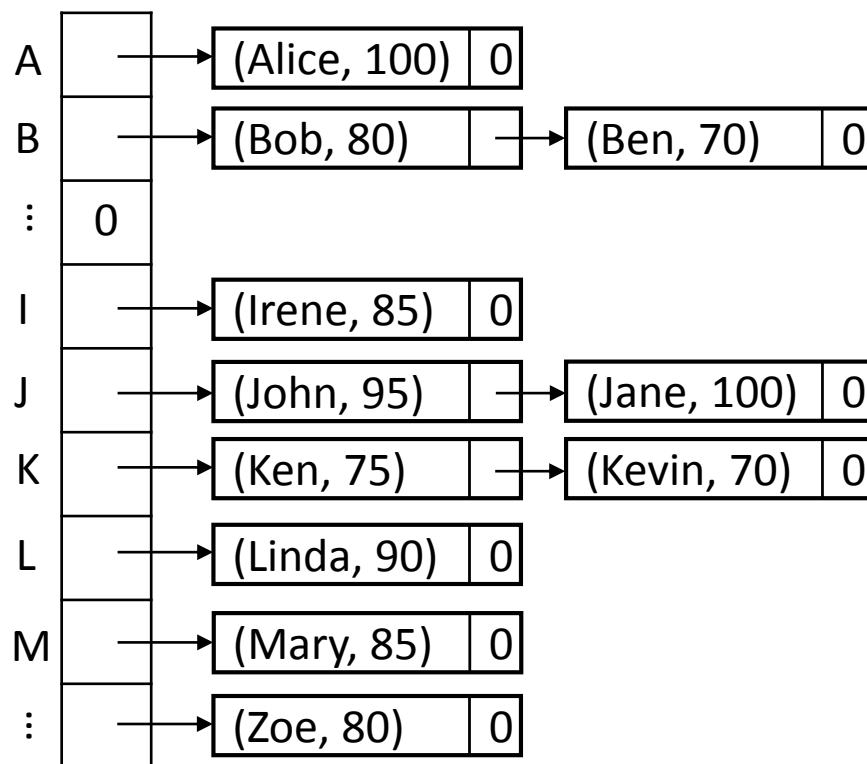


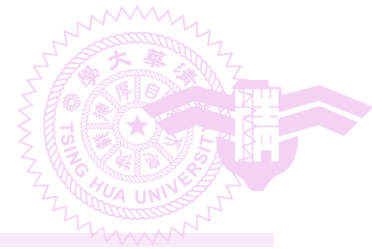
- Consideration:
 - We need to argue the advantages of our hash compared with the commonly used ones



Chain-Based Hash Table

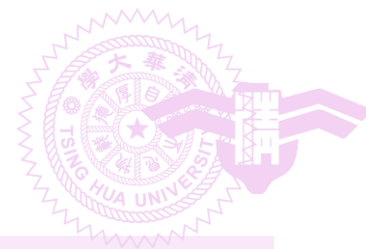
- Each bucket is a chain
 - Chain nodes are typically unordered
 - We typically expect the hash function spreads records uniformly enough
 - Thus each chain does not contain too many nodes
- Linearly traversing a chain is required for inserting, finding, and removing a key





Outline

- 8.1 Introduction
- (8.2 Static hashing)
- (8.3 Dynamic hashing)
- **8.4 Bloom filters**
- Security hash



Bloom Filter Concepts

- Proposed by Burton Howard Bloom in 1970
- A probabilistic data structure
 - For constructing a set and then determining whether some keys is in the set

| | Traditional set data structures, e.g., a BST | Bloom filters |
|---|--|---------------|
| False positive (It could be <i>wrong</i> when it says "Yes") | X | ○ (缺點) |
| False negative (It could be <i>wrong</i> when it says "No") | X | X |
| Easy insertion | ○ | ○ |
| Easy deletion | ○ | X (缺點) |
| Memory space efficiency | Low | High (優點) |

Grocery Shop Example

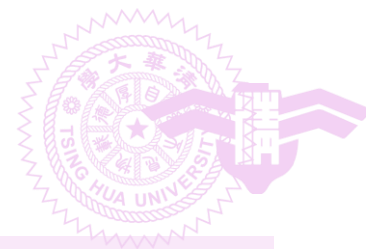
- Suppose we own a grocery shop
- Customers occasionally ask for an item that we are not sure about the availability
 - We spend significant time looking for an item before realizing that the item is unavailable



Grocery Shop Example

- Bloom filter can help
 - Determine the availability of an requested item
 - Some **false positive** are acceptable
 - i.e., the data structure tells that an item is available, but the fact is otherwise
 - No **false negative**
 - We do not want to mistakenly turn down a customer's request





Bloom Filter

- Components
 - A bit vector
 - Multiple hash functions
- Example
 - A table with 26 entries, A ~ Z
 - Three hash functions for a string
 - First character
 - Second character
 - Third character

| | | | |
|---|--|---|--|
| A | | N | |
| B | | O | |
| C | | P | |
| D | | Q | |
| E | | R | |
| F | | S | |
| G | | T | |
| H | | U | |
| I | | V | |
| J | | W | |
| K | | X | |
| L | | Y | |
| M | | Z | |

Bloom Filter

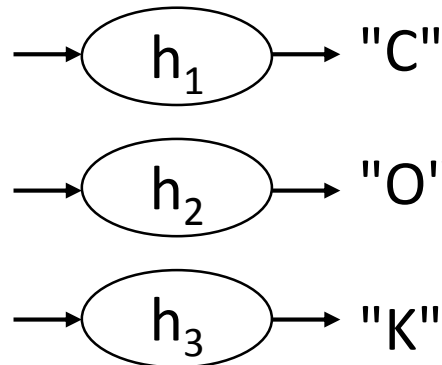
Available items



- Example

- Register string "Coke" into the Bloom filter to indicate that our grocery sells Coke
 - Set the bit vector according to the three hash values, C, O, and K

"Coke"



| | | | |
|---|---|---|---|
| A | | N | |
| B | | O | 1 |
| C | 1 | P | |
| D | | Q | |
| E | | R | |
| F | | S | |
| G | | T | |
| H | | U | |
| I | | V | |
| J | | W | |
| K | 1 | X | |
| L | | Y | |
| M | | Z | |

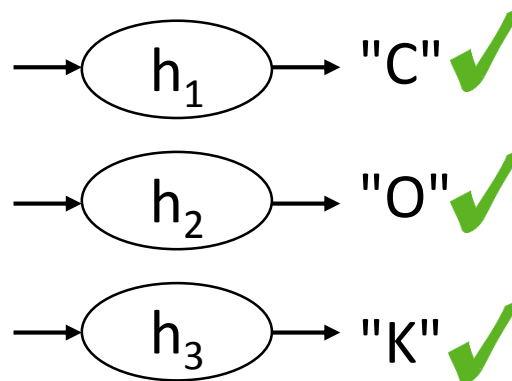
Bloom Filter

Available items



- A simple test
 - If a customer request for "Coke" afterward
 - Bit vector is examined according to the three hash values
 - Bloom filter determines that coke is available because the corresponding bits have been set

"Coke"



| | | | |
|---|---|---|---|
| A | | N | |
| B | | O | 1 |
| C | 1 | P | |
| D | | Q | |
| E | | R | |
| F | | S | |
| G | | T | |
| H | | U | |
| I | | V | |
| J | | W | |
| K | 1 | X | |
| L | | Y | |
| M | | Z | |

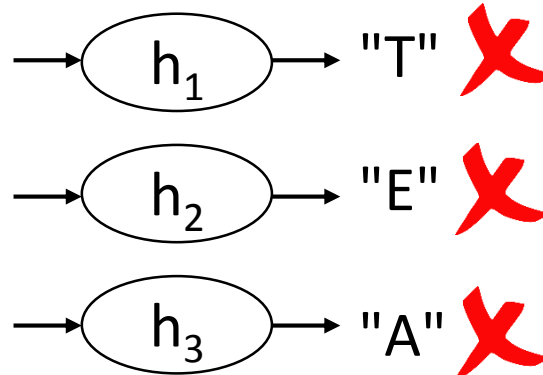
Bloom Filter

Available items



- A simple test
 - If a customer request for "orange juice" afterward
 - Bloom filter determines that orange juice is unavailable because at least one corresponding bit is not set

"Tea"



| | | | |
|---|---|---|---|
| A | | N | |
| B | | O | 1 |
| C | 1 | P | |
| D | | Q | |
| E | | R | |
| F | | S | |
| G | | T | |
| H | | U | |
| I | | V | |
| J | | W | |
| K | 1 | X | |
| L | | Y | |
| M | | Z | |



Bloom Filter

- We register more strings into the Bloom filter



"Fanta" → F A N



"Sprite" → S P R



"Vitali" → V I T

Available items



| | | | |
|---|---|---|---|
| A | 1 | N | 1 |
| B | | O | 1 |
| C | 1 | P | 1 |
| D | | Q | |
| E | | R | 1 |
| F | 1 | S | 1 |
| G | | T | 1 |
| H | | U | |
| I | 1 | V | 1 |
| J | | W | |
| K | 1 | X | |
| L | | Y | |
| M | | Z | |

Bloom Filter

- Test again
 - Bloom filter still works



"Coke" → C O K



"Tea" → T E A



"Fanta" → F A N



Available items



| | | | |
|---|---|---|---|
| A | 1 | N | 1 |
| B | | O | 1 |
| C | 1 | P | 1 |
| D | | Q | |
| E | | R | 1 |
| F | 1 | S | 1 |
| G | | T | 1 |
| H | | U | |
| I | 1 | V | 1 |
| J | | W | |
| K | 1 | X | |
| L | | Y | |
| M | | Z | |



Advantages

- "Coca Cola"
- "Fanta"
- "Sprite"
- "Vitali"
- ⋮

→ Only 26 bits

| | | | |
|---|---|---|---|
| A | 1 | N | 1 |
| B | | O | 1 |
| C | 1 | P | 1 |
| D | | Q | |
| E | | R | 1 |
| F | 1 | S | 1 |
| G | | T | 1 |
| H | | U | |
| I | 1 | V | 1 |
| J | | W | |
| K | 1 | X | |
| L | | Y | |
| M | | Z | |

Disadvantages

- Bloom filter exhibits **false positive**
 - When Bloom filter says "yes", it is not 100% true
 - But, when Bloom filter says "no", it is always true
- "Coffee" is a false positive in our example



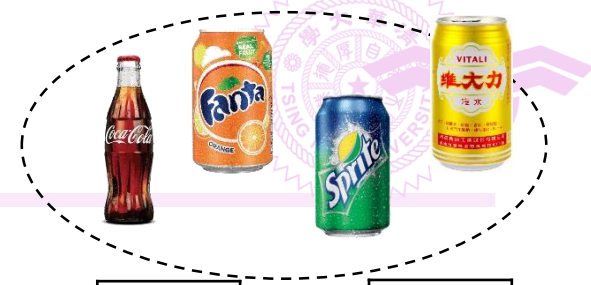
'Coffee' → C O F



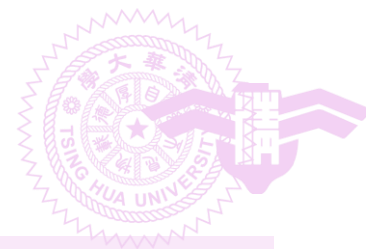
Our grocery does not sell coffee actually!



Available items

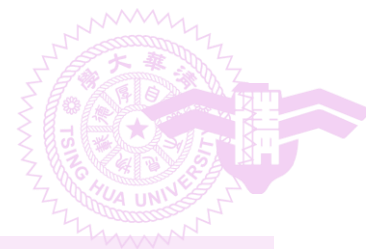


| | | | |
|---|---|---|---|
| A | 1 | N | 1 |
| B | | O | 1 |
| C | 1 | P | 1 |
| D | | Q | |
| E | | R | 1 |
| F | 1 | S | 1 |
| G | | T | 1 |
| H | | U | |
| I | 1 | V | 1 |
| J | | W | |
| K | 1 | X | |
| L | | Y | |
| M | | Z | |



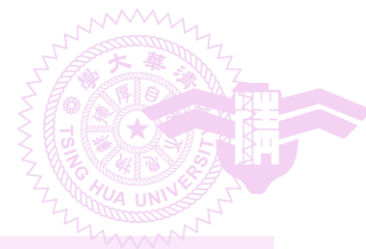
Bloom Filter Analysis

- Key factors of a bloom filter
 - Number of hash functions, k
 - Number of bits in the bit vector, m
 - Number of items expected to be stored, n
 - Uniformity of the hash functions
- False positive analysis
 - Bit vector is set nk times after n items are stored
 - Each time, the probability that a particular bit is set is $(1/m)$
 - Assume true uniformity of hash functions
 - The probability that a bit is set is $(1 - (1 - 1/m)^{nk})$ after n items are stored
 - The probability of a false positive is $(1 - (1 - 1/m)^{nk})^k$
- We can carefully select m , n , and k to achieve our acceptable false positive rate, e.g., 1%



Outline

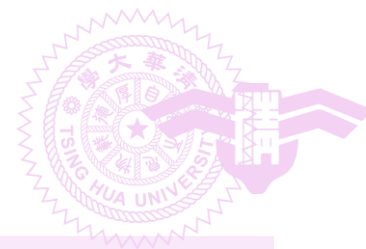
- 8.1 Introduction
- (8.2 Static hashing)
- (8.3 Dynamic hashing)
- 8.4 Bloom filters
- **Security hash**



Security Hash

- Example
 - MD5
 - SHA256
- Usage
 - Password store (能驗證密碼但又可防範洩漏密碼)
 - Digital signature (防止變造)
 - Digital currency
 - ...






Security Hash: MD5

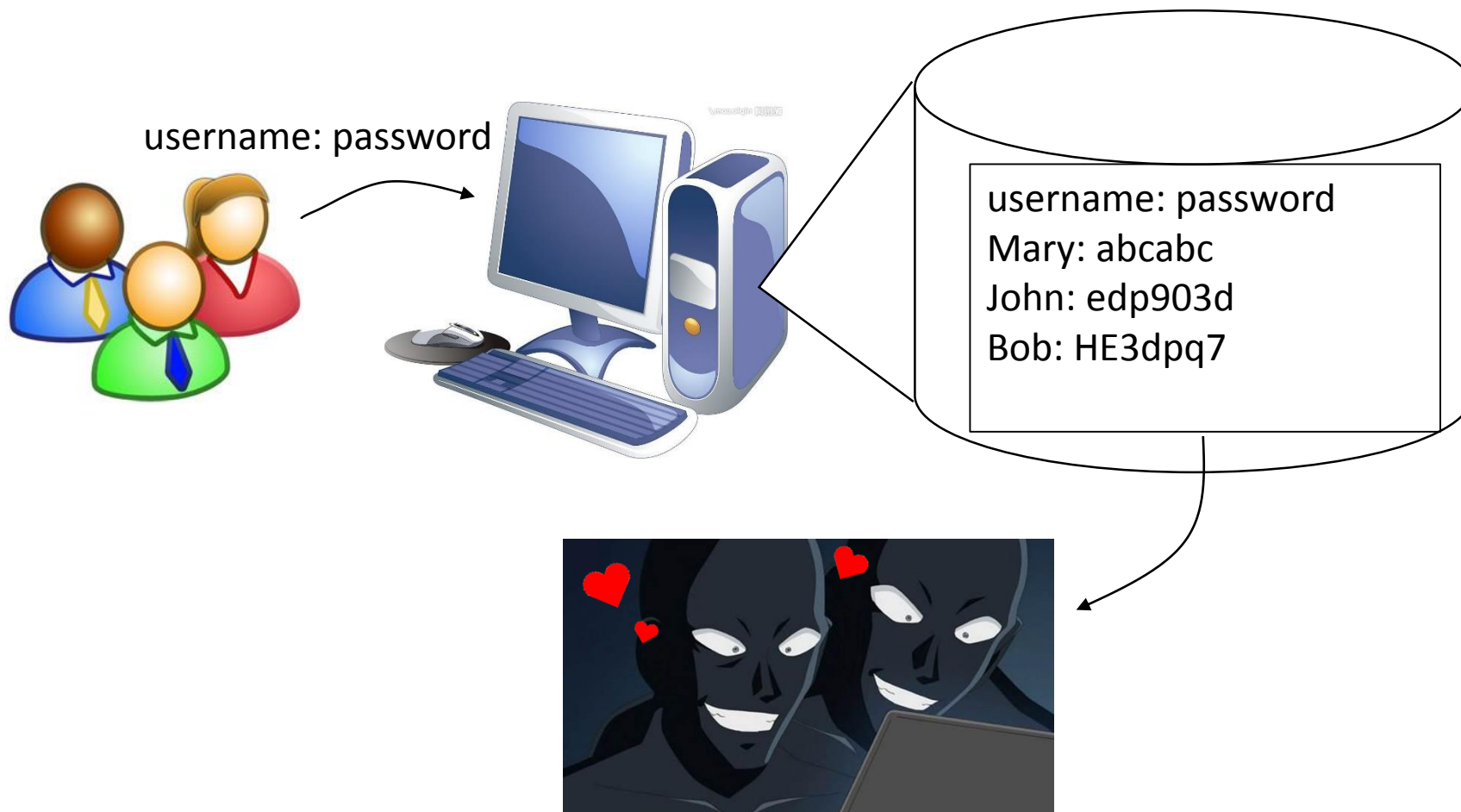
- <https://www.md5hashgenerator.com/>
- Example
 - “NTHU” → 8191af722cfd2890b7a9e986003a6439
 - “NTHU1” → 4c97870289739e75576a4cbeb6222e25
 - “資料結構” → 9aa8415c41bb436d4b6e5618a7be6360



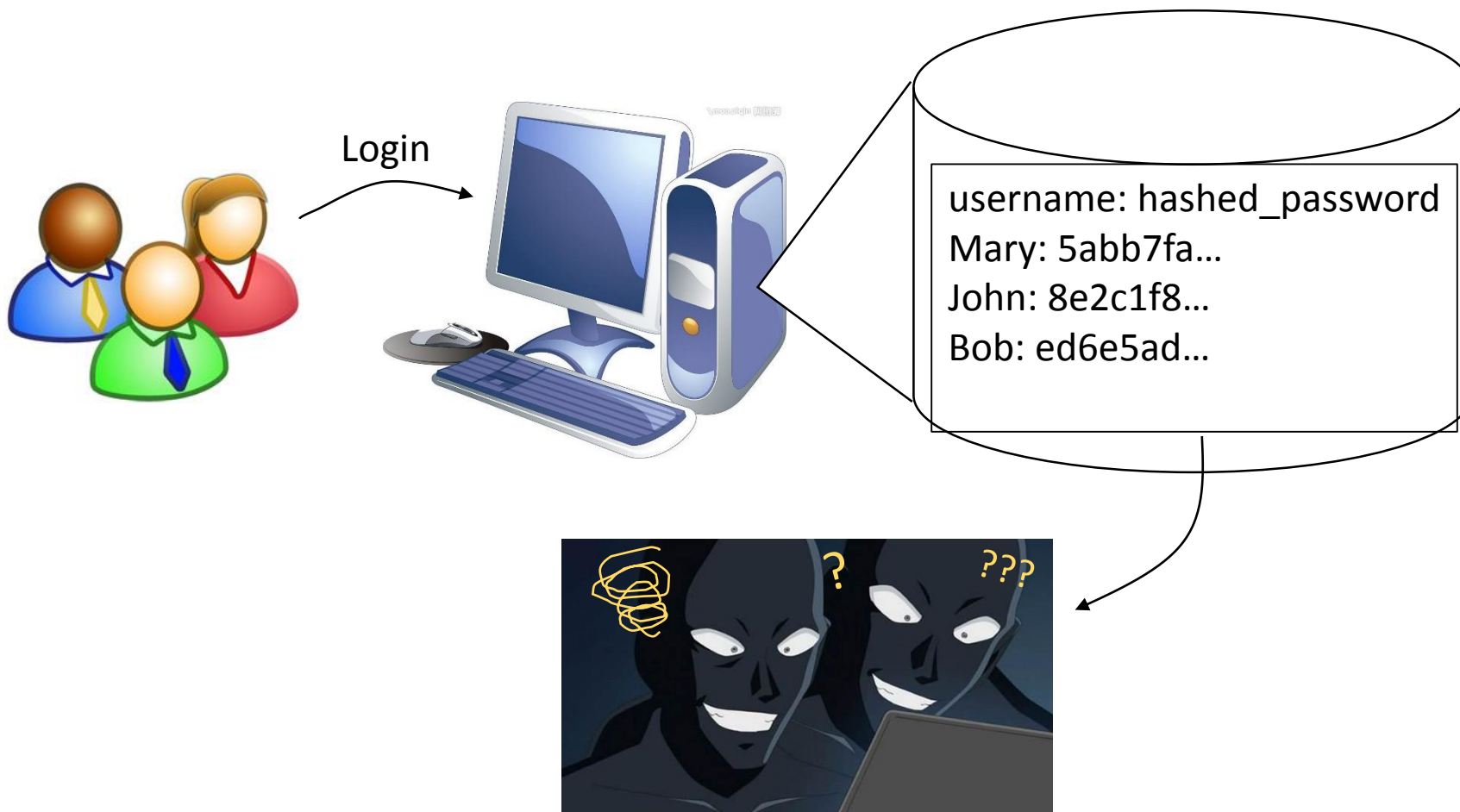
Security Hash

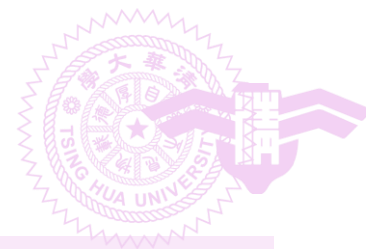
- Important properties
 - Hard to find the inverse
 - “” →
21cee26fd407729a1c740105891e3fca
 - Easy to verify
 - Hard to find another meaningful input that has the same hash

Usage: Password Store



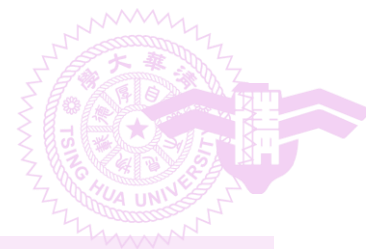
Usage: Password Store





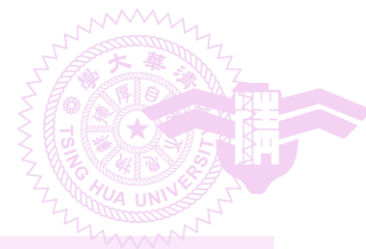
Usage: Digital Signature

- 有獎徵答，請大家解一題資結問題，前3快解出問題的給獎品!
 - 狀況一
 - 有人搶先說“我解開了”再慢慢解題
 - 狀況二
 - 第一名公布了答案，二三名都說“我也是這麼想的”
 - 狀況三
 - 大家統一把答案交給某個裁判再檢查是誰最快回答出正確答案。但裁判也有可能洩漏答案?
- Security hash可以幫忙解決這個問題



Usage: Digital Currency

- 數位貨幣(例如比特幣)
 - 去中心化
 - 由貨幣使用者記帳，而非由政府或銀行記帳
 - 記帳的使用者會得到一些報酬
- 問題
 - 怎麼決定記帳權歸屬，且不會被特定使用者把持？
 - 怎麼避免帳本被竄改？



Usage: Digital Currency

- 方法
 - 要求使用者解一hash inverse問題
 - 請把x取代成數字，使" A轉給B十元(xxxxxxxx)"的MD5前N碼是0
 - A轉給B十元(00000000) → **4299bf747...**
 - A轉給B十元(00000001) → **2c81a694d...**
 - A轉給B十元(00000002) → **551741357...**
 - 求解動作稱挖礦，困難度被控制在約10分鐘才能解答一次。
 - 挖礦有利可圖，因此很多人參與
 - 想要單方把持記帳權，必須投入非常多電腦，不敷成本
 - 後帳本與前帳本相關(區塊鏈)
 - 想從竄改已經成形的帳本，必須投入非常多電腦在短時間內重算竄改位置之後的hash inverse問題，不敷成本