# Hashing

---

## Hashing

- Start with an *array* that holds the hash table.
- Use a *hash function* to take a key and map it to some index in the array.
- If the desired record is in the location given by the index, then we are finished; otherwise we must use some method to resolve the *collision* that may have occurred between two records wanting to go to the same location.

# Today's Math

$$U = 1 + 2x + 3x^2 \ldots\ldots + mx^{m-1} = \sum_{k=1}^{m} k \cdot x^{k-1}$$

- We know:

$$\frac{x^{m+1} - 1}{x - 1} = 1 + x + x^2 \ldots\ldots + x^m$$

$$\frac{(x-1)(m+1)x^m - (x^{m+1} - 1)}{(x-1)^2} = 1 + 2x + 3x^2 \ldots\ldots + mx^{m-1}$$

$$if\, |x| < 1$$

$$\sum_{k=1}^{\infty} kx^{k-1} = \lim_{m \to \infty} \frac{(x-1)(m+1)x^m - (x^{m+1} - 1)}{(x-1)^2} = \frac{1}{(x-1)^2}$$

---

# Choice of Hash Function

- Quick to Compute
- Randomization

- Truncation
    - pick first second and fifth numbers
    - example n=62538194, h=394
- Folding
    - use all
    - 62538194 maps to 625+381+94=1100
- Modular Arithmetic
    - n mod HASHSIZE. HASHSIZE is some prime
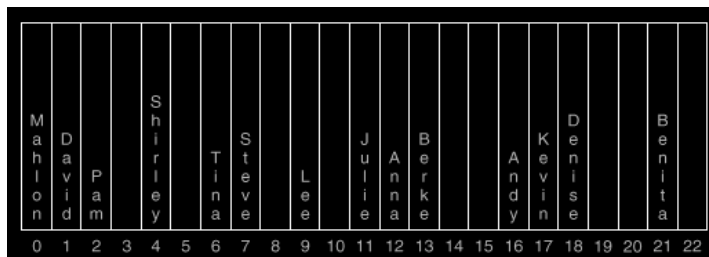    - 47 mod 7 =5

## Collision Resolution by Open Addressing

- Linear Probing
  - *Linear probing* **starts with the hash address and searches sequentially for the target key or an empty position. The array should be considered circular, so that when the last location is reached, the search proceeds to the first location of the array.**

---

## Collision Resolution by Open Addressing

- Clustering
  - *Linear probing* **starts with the hash address and searches sequentially for the target key or an empty position. The array should be considered circular, so that when the last location is reached, the search proceeds to the first location of the array.**
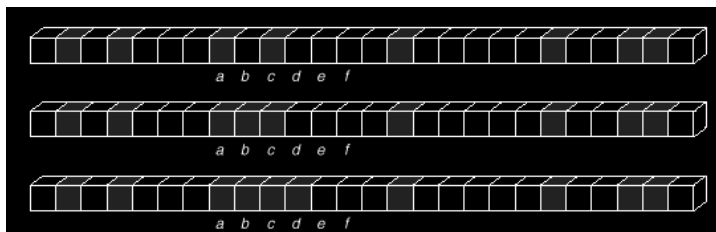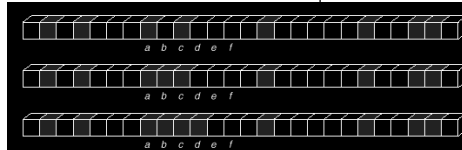
# Collision Resolution by Open Addressing

- Quadratic Probing
  - If there is a collision at hash address *h*, *quadratic probing* goes to locations *h*+1, *h*+4, *h*+9, that is, at locations $h+i^2$ (mod hashsize) for $I = 1,2...$

> **Quadratic Probing only searches half of the locations.**

- Other Probing Methods
  - Key dependent probing
  - Random Probing

---

# Key Deletion

- Simple Deletion from Hash Table:
  - 11=h(Julie); 12=h(Anna);11=h(Berke)
  - Now we delete Anne.
  - Can we find **Berke**?
- Solution?

| M a h l o n | D a v i d | P a m | | S h i r l e y | | T i n a | S t e v e | | L e e | | J u l i e | A n n a | B e r k e | | | A n d y | K e v i n | D e n i s e | | | B e n i t a | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 |

*4*

# Collision Resolution by Chaining

---

# Pros and Cons of Chaining

- Simple and efficient collision handling
- No Overflow
- Easy Deletion
- Space saving if records are large. The size of static hash table is still small. Only the chain grows.
- Cons: Extra space in links. The relative waste increases if records are small.

## Birthday Surprise

- With randomly chosen people in a room, what is the probability that no two have the same birthday?
    - The probability that the second person has no birthday collision is 364/365
    - The probability that the second person has no birthday collision is 363/365
    - The probability that mth person has a different birthday is (365-m+1)/365
    - The probability that all m persons have separate birthday:
      $$\frac{364}{365} \times \frac{363}{365} \times .....\times \frac{365-m+1}{365}$$
    - This becomes less than .5 when m>23!

# Analysis of Hashing

12

## Definitions

- What is the cost factor?
  - A *probe* is one comparison of a key with the target.
- Load factor
  - The *load factor* of the table is $\lambda = n/t$, where $n$ positions are occupied out of a total of $t$ positions in the table.

## Analysis of Chaining

- Unsuccessful retrieval
  - a chain have to be searched until end.
  - Average chain is $\lambda = n/t$

- Successful retrieval
  - (n-1) mismatched keys and 1 matching key.
  - Average mismatch keys per chain (n-1)/t
  - Average probe (n-1)/2.t +1 $\approx$ 1+ $\lambda/2$

# Analysis of Open Addressing (random probe)

- Unsuccessful Probe:
  - An unsuccessful search terminates when it encounters an empty space.
  - Probability that the first probe hits a full cell= $\lambda$
  - Probability that the first probe hits an empty cell =$(1-\lambda)$
  - Probability for exact two probe and it termites is = $\lambda.(1-\lambda)$
  - Probability for exact k probe =$\lambda^{(k-1)}(1-\lambda)$
  - The expected number of probe:

$$U(\lambda) = \sum_{k=1}^{\infty} k \cdot \lambda^{k-1}(1-\lambda) = ??$$

$$= \frac{1}{1-\lambda}$$

---

# Analysis of Open Addressing (random probe)

- Successful Probe:
  - A successful probe will be equal to the number of unsuccessful search made before inserting the entry, plus one.
  - The table is initially empty with load=0, and it grows.
  - The average number of search in a successful search is:

$$S(\lambda) = \frac{1}{\lambda} \int_{0}^{\lambda} U(\mu)d\mu = \frac{1}{\lambda} \ln \frac{1}{1-\lambda}$$

## Analysis of Open Addressing (linear probe)

- A little more complex analysis since, successive probes are dependant.

Retrieval from a hash table with open addressing, linear probing, and load factor $\lambda$ requires approximately

$$\frac{1}{2}\left(1 + \frac{1}{1-\lambda}\right) \quad \text{and} \quad \frac{1}{2}\left(1 + \frac{1}{(1-\lambda)^2}\right)$$

probes in the successful case and in the unsuccessful case, respectively.

---

| Load factor | 0.10 | 0.50 | 0.80 | 0.90 | 0.99 | 2.00 |
|---|---|---|---|---|---|---|
| *Successful search, expected number of probes:* | | | | | | |
| Chaining | 1.05 | 1.25 | 1.40 | 1.45 | 1.50 | 2.00 |
| Open, Random probes | 1.05 | 1.4 | 2.0 | 2.6 | 4.6 | — |
| Open, Linear probes | 1.06 | 1.5 | 3.0 | 5.5 | 50.5 | — |
| *Unsuccessful search, expected number of probes:* | | | | | | |
| Chaining | 0.10 | 0.50 | 0.80 | 0.90 | 0.99 | 2.00 |
| Open, Random probes | 1.1 | 2.0 | 5.0 | 10.0 | 100. | — |
| Open, Linear probes | 1.12 | 2.5 | 13. | 50. | 5000. | — |

**Theoretical Comparisons**

| Load factor | 0.1 | 0.5 | 0.8 | 0.9 | 0.99 | 2.0 |
|---|---|---|---|---|---|---|
| *Successful search, average number of probes:* | | | | | | |
| Chaining | 1.04 | 1.2 | 1.4 | 1.4 | 1.5 | 2.0 |
| Open, Quadratic probes | 1.04 | 1.5 | 2.1 | 2.7 | 5.2 | — |
| Open, Linear probes | 1.05 | 1.6 | 3.4 | 6.2 | 21.3 | — |
| *Unsuccessful search, average number of probes:* | | | | | | |
| Chaining | 0.10 | 0.50 | 0.80 | 0.90 | 0.99 | 2.00 |
| Open, Quadratic probes | 1.13 | 2.2 | 5.2 | 11.9 | 126. | — |
| Open, Linear probes | 1.13 | 2.7 | 15.4 | 59.8 | 430. | — |

**Empirical Comparisons**

# Comments

- Chaining consistently requires fewer probing than open addressing.
- Which method to use when unsuccessful search is more common?
- If most cases are successful, and the table is not nearly full simpler method of linear probing is not significantly slower than other complex methods.

| Load factor | 0.10 | 0.50 | 0.80 | 0.90 | 0.99 | 2.00 |
|---|---|---|---|---|---|---|
| *Successful search, expected number of probes:* | | | | | | |
| *Chaining* | 1.05 | 1.25 | 1.40 | 1.45 | 1.50 | 2.00 |
| *Open, Random probes* | 1.05 | 1.4 | 2.0 | 2.6 | 4.6 | — |
| *Open, Linear probes* | 1.06 | 1.5 | 3.0 | 5.5 | 50.5 | — |
| | | | | | | |
| *Unsuccessful search, expected number of probes:* | | | | | | |
| *Chaining* | 0.10 | 0.50 | 0.80 | 0.90 | 0.99 | 2.00 |
| *Open, Random probes* | 1.1 | 2.0 | 5.0 | 10.0 | 100. | — |
| *Open, Linear probes* | 1.12 | 2.5 | 13. | 50. | 5000. | — |

---

# Comments (contd..)

- In Hashing based IR, the retrieval time is dependant on load factor not on the number of elements in the list.
  - 20,000 keys in a hash table of 40,000 is same as 20 keys in a list of 40!
- The key to performance is the hash function
  - how quickly it can be evaluated
  - how well it spread the data.

## Game of Life Revisited..

- In version 2 we solved the problem of sparse computation.
- How about space complexity?
  - Perhaps hashing can help!
- For each cell we need to keep:
  - status (live or dead)
  - neighbor count
  - x and y
- Open Addressing or Chaining?
  - Large hash table vs. 25% pointer overhead.
- 4 way linked list.
  - Each node must be a member of four lists maylive, maydie, newlive and newdie.

# MIDTERM
# Review

22

# Midterm

- 4 Questions Total:
  - 1 True-False
  - 1 Overall Concept
  - 1 Searching & Sorting
  - 1 Hashing
- Open Book 60 min.