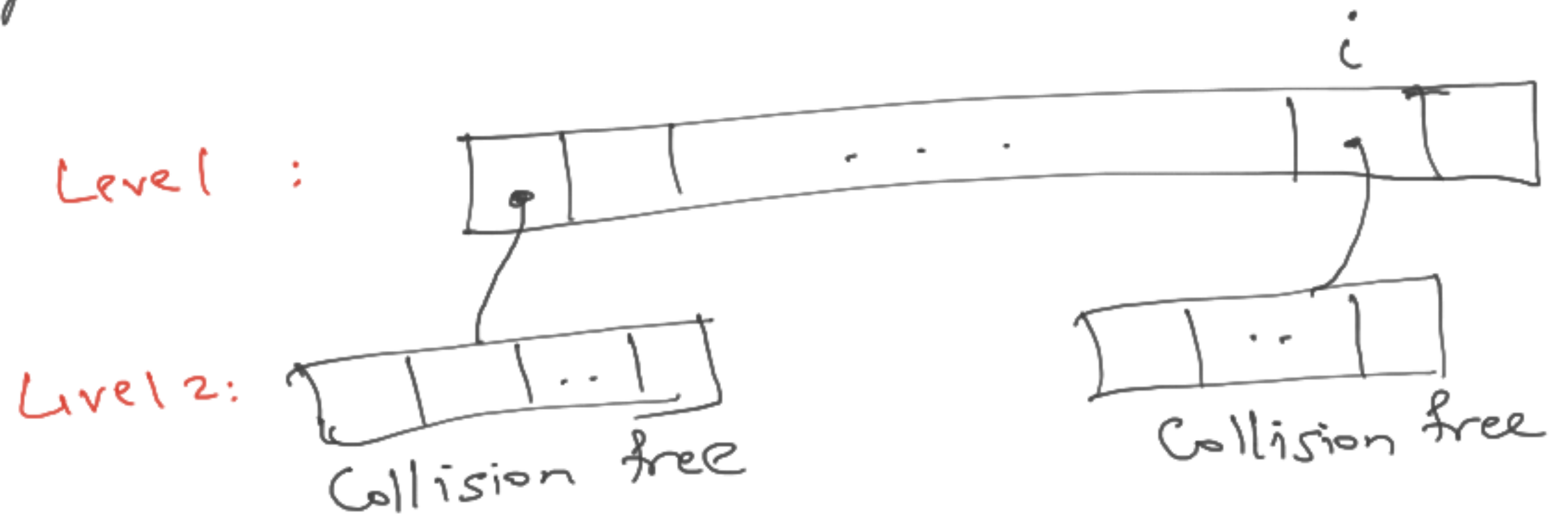


Hashing 2

Two-level hashing ("perfect hashing")



Let $C(i)$ = num. of elements that gets mapped to location i
(in the first level)

∴ Level 2 HT is $C(i)^2 \rightarrow$ Collision free for location i

Total table space for Level 2

$$= \sum_{i=1}^M c(i)^2$$

We know

$$E[C] = \binom{N}{2} \frac{1}{M}$$

$$E \left[\sum_{i=1}^M \binom{c(i)}{2} \right] = \binom{N}{2} \frac{1}{M}$$

$$E \left[\sum c(i)^2 - \sum c(i) \right] = O(N)$$

$$\Rightarrow \underline{E \left[\sum c(i)^2 \right] = O(N)}$$

Collision-free and $O(N)$ table space!

(since $M = O(N)$)

More stronger properties:

k-wise independent hash functions

Defn: $H: U \rightarrow [M]$ k-wise indep. if
for any k distinct keys x_1, \dots, x_k and values

$\alpha_1, \alpha_2, \dots, \alpha_k$

$$P(h(x_1) = \alpha_1 \cap h(x_2) = \alpha_2 \cap \dots \cap h(x_k) = \alpha_k) \leq \frac{1}{M^k}$$

$\rightarrow k=2$, pair-wise indep.

Properties: H is k -wise indep for $k \geq 2$. Then

1. H is also $(k-1)$ -wise indep.

2. For any $x \in U$ and $a \in [M]$, $P[h(x) = a] \leq \frac{1}{M}$

3. H is universal

Pairwise indep \rightarrow stronger! vs. universal?

$$h(x) = Ax$$

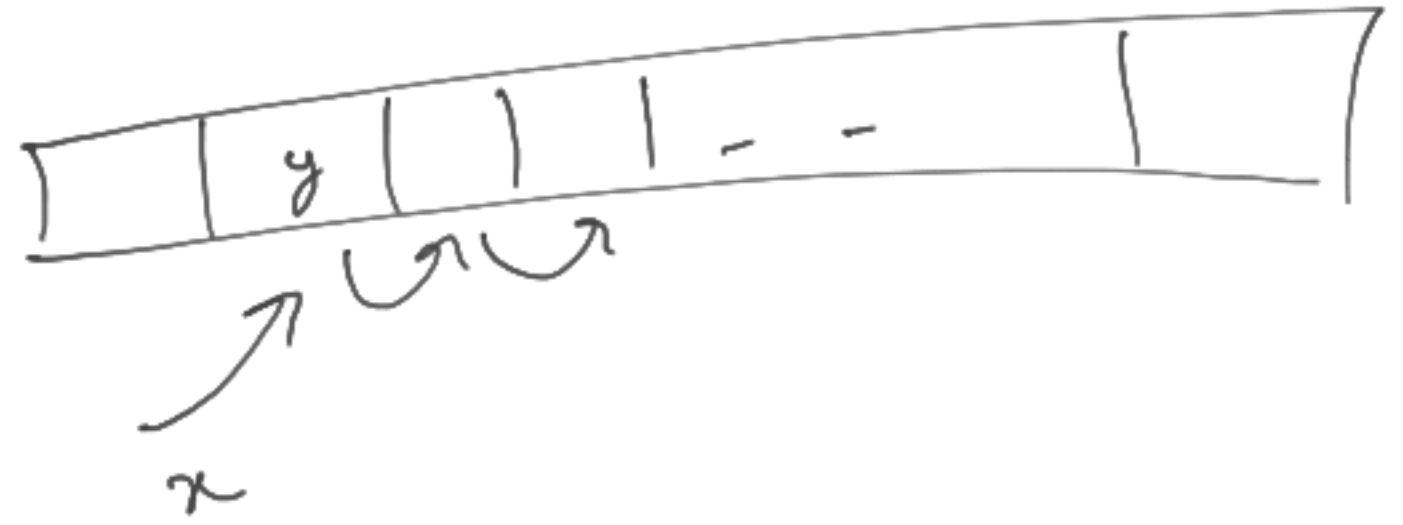
Consider $x = 0$

$$h(0) = 0$$

To turn this into pairwise indep : $h(x) = Ax + b$
 \uparrow
 Uniform random binary vec.

Open addressing

- Single array
- No separate D.S.
- Linear probing
 - use step-size



Quadratic

Cuckoo Hashing

Two tables T_1 & T_2

both of size $M = O(N)$

Two hash functions h_1 & $h_2 \in H$

assume fully random

($O(\log N)$ -wise indep suffices)

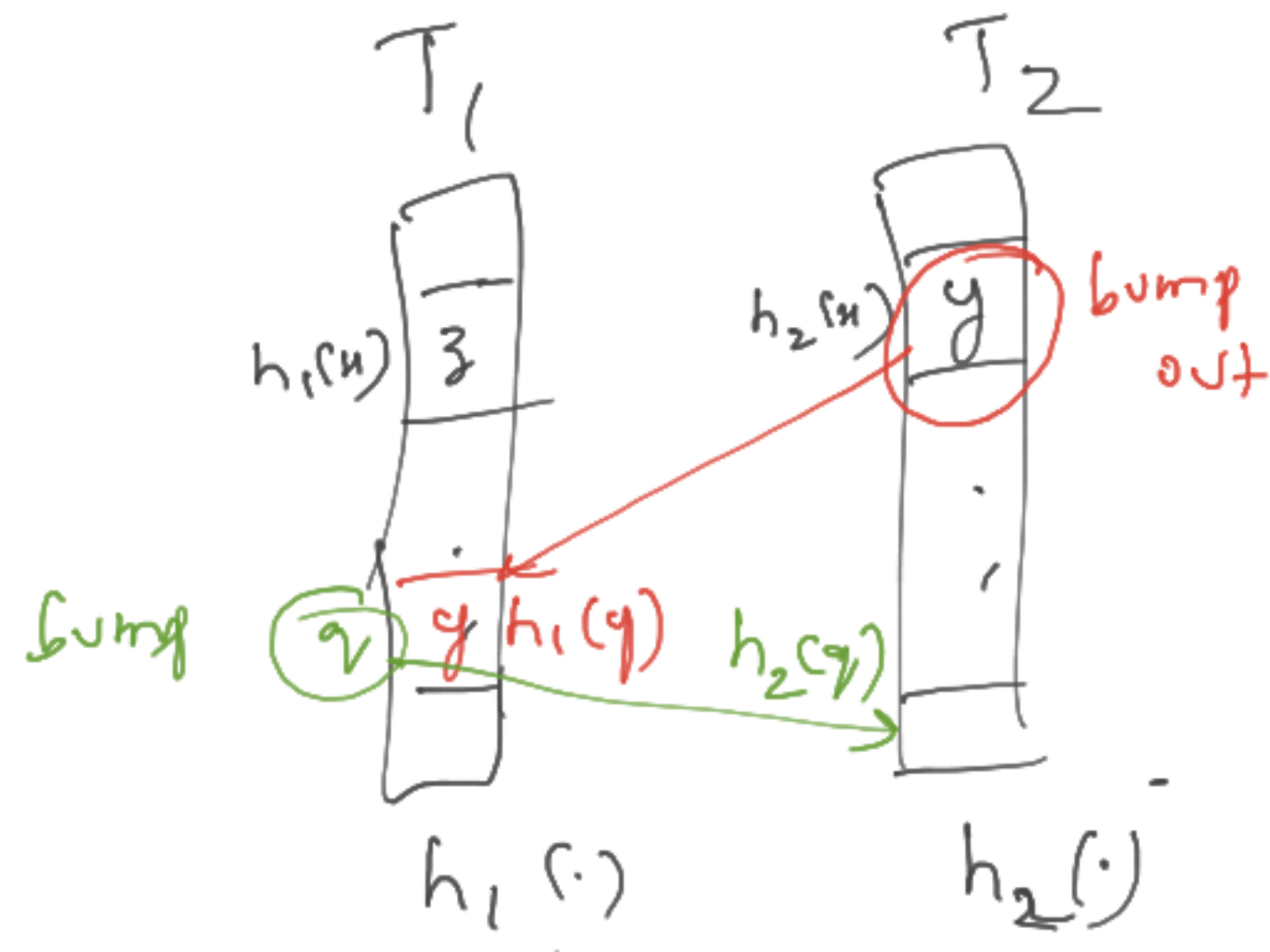
Insertion:

x goes into either $T_1[h_1(x)]$ or $T_2[h_2(x)]$

- stop when no more bumps
- or if more than $C \cdot \log N$ bumps & rehash

Query:

$\Theta(1)$ only 2 locations.



Thm: The expected time to perform insert is $O(1)$
if $M \geq 4N$

Proof Sketch:

"Cuckoo graph" G .

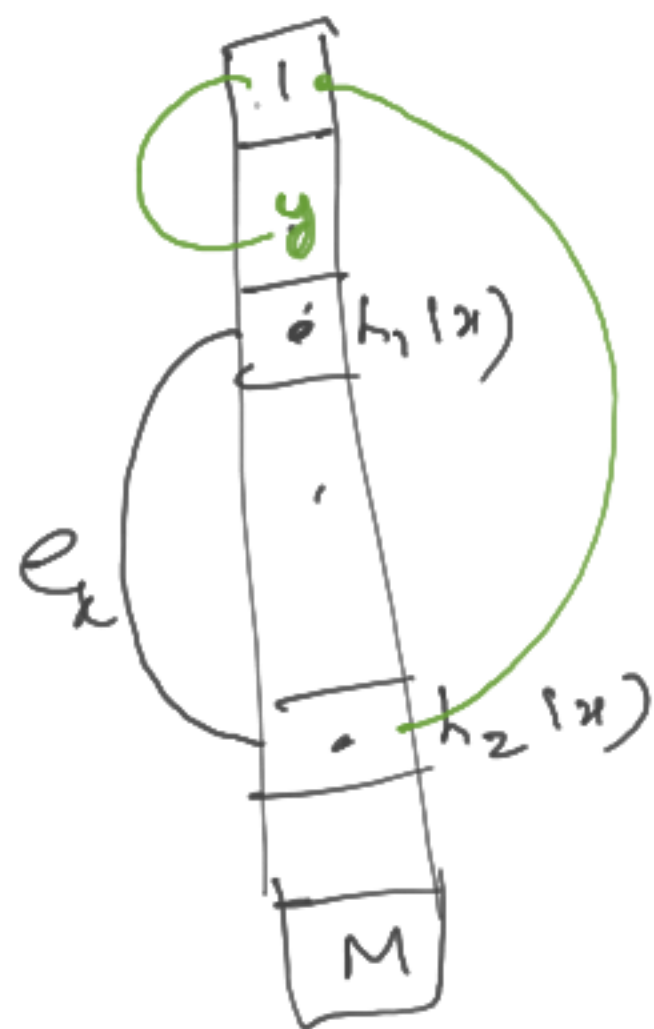
M vertices:

= hash table locations.

Edges: correspond to the items to be inserted.

$\forall x \in S, e_x = (h_1(x), h_2(x))$
→ for all

Cuckoo graph



"y" can get bumped out when inserting a new element x
if y falls in a path in the cuckoo graph
starting from $h_1(x)$ or $h_2(x)$

Defining: Bucket of x
 $B(x) =$ set of nodes of G reachable from
 $h_1(x)$ or $h_2(x)$
= 'connected component' of G with edge e_x

$E[\text{insertion time of } x] = E[|B(x)|]$
size of $B(x)$

To show: $E[|B(x)|] \leq O(1)$

$$E[|B(x)|] = \sum_{\substack{y \in S \\ y \neq x}} P[e_y \in B(x)]$$

$$\leq N P[e_y \in B(x)]$$

Union bound

Sufficient to show

$$P[e_y \in B(x)] \leq O\left(\frac{1}{M}\right)$$

$M \geq 4N$.

Lemma: For any $i, j \in [M]$

P [there exists a path of length l between i and j in the cuckoo graph] $\leq \frac{1}{2^l M}$

Proof: (exercise)

For $l=1$. . .

Then use induction .

To show: $P[e_y \in B(x)] \leq \Theta\left(\frac{1}{M}\right)$

Ans: Using the lemma

$$P(e_y \in B(x)) \leq \sum_{l \geq 1} \frac{1}{2^l M}$$
$$= O\left(\frac{1}{M}\right)$$

(from
prev lemma)

$$M \geq 4N$$

Application : Bloom Filter

Membership query

Room for mistakes :

Only false positives

but no false negatives

Useful for 'filter operations' - typically elements not in the set

Space efficient data structure for approximate membership queries

• Array T of M bits

• k hash function $h_1, h_2, \dots, h_k: U \rightarrow [M]$
(Assume completely random)

Adding a key:

$x \in S$

set bits

$T[h_1(x)], T[h_2(x)], \dots, T[h_k(x)]$

to 1

Membership query: check locations

$x: T[h_1(x)] \dots T[h_k(x)]$

Let $p = \text{prob. that a bit in } T \text{ is not set}$

$$p = \left(1 - \frac{1}{M}\right)^{kN} = \left(1 - \frac{1}{M}\right)^{M \cdot \frac{kN}{M}}$$

$$= \left(e^{-1}\right)^{\frac{kN}{M}}$$

$$= \underline{\underline{e^{-\frac{kN}{M}}}}$$

$$\left(1 - \frac{1}{2}\right)^x \rightarrow \frac{1}{e}$$

Prob. of false positive = all k bits set

$$= (1 - p)^k = \left(1 - e^{-\frac{kN}{M}}\right)^k$$

$$\frac{d}{dk} () = 0 \quad (\text{exercise})$$

False positive prob minimized

$$\underline{\underline{k = \frac{M}{N} \ln 2}}$$

↓ ϵ detector prob. of false positive.

$$\epsilon = \left(\frac{1}{2}\right)^{\frac{M}{N} \ln 2}$$

$$\underline{\underline{M = 1.44 \log(1/\epsilon)}}$$

For 1% false positive.

$$M \leq 10 N \text{ bits}$$

$$k = 7$$

$$\underline{\underline{N * \log |U|}}$$

