# Hashing:

We briefly encountered a similar problem when dealing with balls and bins.

If $n$ balls were placed in $n$ bins uniformly at random, no bin contains more than $\frac{3\ln(n)}{\ln\ln(n)}$ balls w.p $\geq 1 - \frac{1}{n}$.

Number of empty bins was concentrated around its expectation of $\frac{n}{e}$.

1. Space + time to be $O(n)$
2. Additions $\sim O(1)$ time
3. Query $\sim O(1)$ time.

# Perfect hashing: ✓

Given a set $S$ of $n$ keys from the universe $U$, build a lookup table of size $O(n)$ s.t a membership query can be answered in $O(1)$ time.

⤷ we call it "perfect hashing for $S$".

Defn: A set of hash functions $H$ is called a weak universal family if for all $x, y \in U$, $x \neq y$.     $H: U \rightarrow [m]$.

$$\Pr_{h \in H}\left[h(x) = h(y)\right] = \frac{O(1)}{m}.$$

Fix an element $x$.

Expected chain length $= \underset{h \sim H}{E}\left[\# \text{ of } y \text{ st } \underset{x \neq y}{h(x) = h(y)}\right] + 1$

We want $1 + \frac{n-1}{m} = O(1)$.

$m$ and $n \cdot m \simeq O(n)$.

$= 1 + \underset{\substack{y \\ h \sim H}}{\sum} \Pr\left[h(x) = h(y)\right]$

$= 1 + (n-1) \cdot O(\frac{1}{m})$

$\leq 1 + \frac{1}{m} \cdot (n-1)$

# Fredman-Komlos-Szemeredi hashing:

→ Expected $O(n)$ time
→ Worst case $O(n)$ space  } [2001].
→ $O(1)$ worst case query time.

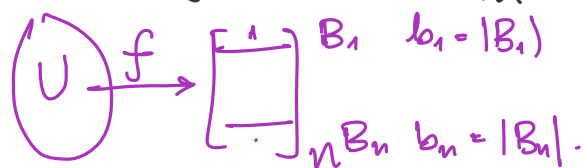Suppose $m = \Omega(n^2)$; Then if we pick $h$ randomly from $H$

Then
$$\mathbb{E}\left[\text{no. of collisions}\right] = \sum_{\substack{x,y \in S \\ x \neq y}} \Pr\left[h(x) = h(y)\right]$$

$$= \binom{n}{2} \cdot \frac{O(1)}{m}$$

$$\leq \underline{O(1)} \cdot \checkmark$$

⇒ $\underline{O(1)}$ trials needed to obtain collision free hashing.

Suppose $m \sim n$ Then

$$\mathbb{E}\left[\text{no. of collisions}\right] = \binom{n}{2} \cdot \frac{O(1)}{m} = \underline{O(n)}.$$


$B_1 \quad b_1 = |B_1|$
$_n B_n \quad b_n = |B_n|$

$m \sim n$

$U = \{0, \ldots, |U|-1\}$

FKS gave a 2-step hashing algo.

· Find a hashing fn $f: U \to [n]$ that partitions
S into buckets $B_1, \ldots, B_n$. ← Obtain this from weak univ
$|B_i| = b_i \Rightarrow \sum b_i = n.$   class of hash functions.

· For each bucket $B_i$: find a function $g_i : U \to [O(b_i^2)]$

Say $|B_i| = b_i$.

For each element $x \in S$, let $C_x$ be the no. of its collisions

$\checkmark\ \underset{h}{\mathbb{E}}[C_x] = \frac{|S|-1}{n} = \frac{n-1}{n} < 1. \quad \checkmark$

$\underset{y \neq x}{\sum} 1 \cdot \underset{h}{Pr}[y \neq x : h(x) = h(y)] \le \frac{n-1}{n} < 1$

1. Select a random function from H.

2. Compute a hash table with chaining, so insertion takes $O(1)$ time

3. Compute an auxillary array $B_2$ s.t $|B_2(i)| = O(b_i^2)$

4. If $\sum_{i=1}^{n} b_i^2 > \beta n$ then go to step 1. Else "record" f.

$\sum b_i^2 \le \beta n.$

Let $t$ be the no. of iterations. We want to bound the expected no. of iterations.

Claim: If $\beta > 4$, $\mathbb{E}[t] \le 2$.

Proof: Total no. of collisions ($C_S$) is as follows

$C_S = \sum_{i=1}^{n} |\{(x,y) \mid x, y \in B_i , x \neq y\}|$

$= \sum_{i=1}^{n} b_i (b_i - 1)$

$\sum_{i=1}^{n} b_i = n.$

$= \sum_{i=1}^{n} (b_i^2 - b_i) = \left(\sum_{i=1}^{n} b_i^2\right) - n$

$< 1$

$$\mathbb{E}_h[C_S] = \mathbb{E}_h\left[\sum_{i=1}^{n} b_i^2\right] - n \quad\longrightarrow\quad \mathbb{E}[C_S] = \sum_{x \in S} \underline{\mathbb{E}[C_x]}$$

We know that $\mathbb{E}[C_S] < n$.

$$\leq \sum_{x \in S} 1 = n.$$

$$\Rightarrow \mathbb{E}\left[\sum b_i^2\right] < 2n.$$

Using Markov's ineq: $\Pr\left[\sum_{i=1}^{n} b_i^2 > 4n\right] < \frac{1}{2}$.

$b_i \rightsquigarrow O(b_i^2)$.

$$\overset{\beta}{\underset{\downarrow}{}}$$

w.p $\geq \frac{1}{2} \quad \sum \underline{b_i^2 \leq 4n}$

$\begin{cases} \bullet \text{ Select } g_i : \underline{U} \rightarrow [\alpha b_i^2] \text{ from } H. \\ \bullet \text{ If for some } \underline{x_i \in B_i} \text{ ; there is a collision, pick a new } \\ \quad \underline{g_i}. \end{cases}$

<span style="color:red">Claim: If $\alpha \geq 2$ then $\mathbb{E}[t_i] = O(b_i^2)$.</span>

Proof: $C_x = $ no. of collisions of $x$ in $B_i$

$$\underset{g_i}{\mathbb{E}[C_x]} < \frac{b_i}{\alpha b_i^2} = \frac{1}{\alpha b_i}. \qquad \underset{g_i \sim H}{\Pr\left[g_i(x) = g_i(y)\right]} \leq \frac{1}{\alpha b_i^2}$$

$$\underset{x \neq y}{}$$

From Markov's ineq.

$$\underset{g_i}{\Pr[C_x \geq 1]} \leq \underset{g_i}{\mathbb{E}[C_x]} < \frac{1}{\alpha b_i}. \qquad\qquad B_i$$

Over all elements in $\underline{B_i}$

$$\underset{g_i}{\Pr[\exists x \text{ st } C_x \geq 1]} \leq \sum_{x} \underset{g_i}{\Pr[C_x \geq 1]} < b_i \cdot \frac{1}{\alpha b_i} = \underline{\frac{1}{\alpha}}.$$

If $\alpha = 2$ then w.p $\geq \frac{1}{2}$, no collisions happen.

$\Rightarrow$ At most 2 trials to get $g_i$ w/ no collisions.
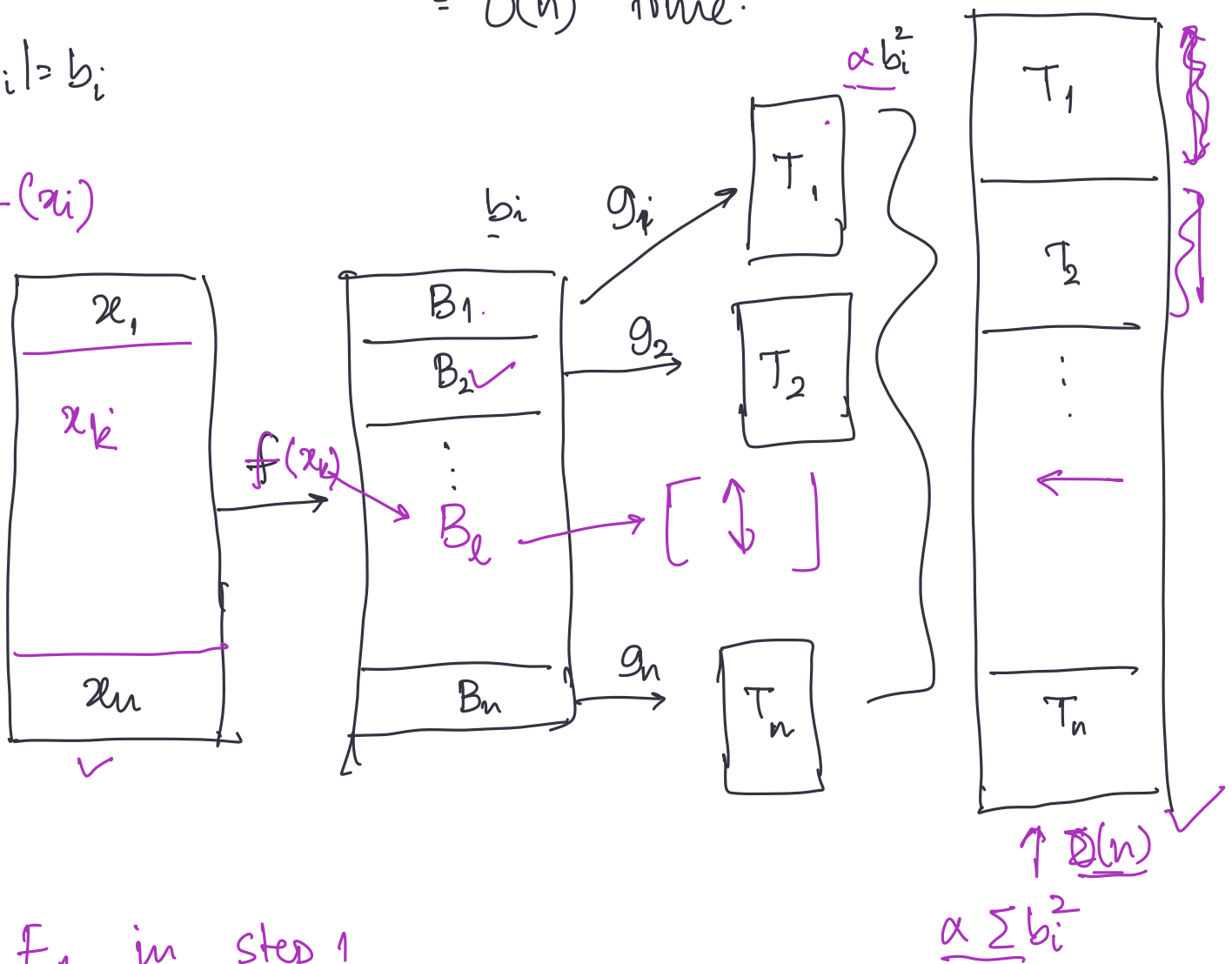
$\Rightarrow$ Expected time $= O(b_i^2)$.
for this step

$\alpha b_i^2$

Total expected running time: $O(n) + \sum O(b_i^2) + O(n)$

preprocessing

$= O(n)$ time.

$|B_i| = b_i$

$f(x_i)$

$\alpha b_i^2$

$b_i \quad g_i$

$T_1$

$x_1$

$x_k$

$f(x_k)$

$B_1$

$B_2$ ✓

$g_2$

$T_1$

$T_2$

$T_2$

$\vdots$

$B_\ell \longrightarrow [\updownarrow]$

$x_n$

$B_n$

$g_n$

$T_n$

$T_n$

$\uparrow O(n)$

$\alpha \sum b_i^2$

$E_1$ in step 1

$E_2$ in step 2.

Total error $\leq E_1 + E_2 < 1$.

$1$