

Lecture 14 DICTIONARIES

(HASHING WITH CHAINING)

① Universal Hashing

[Carter, Wegman '79]

We want to use a hash-function, to store ~~the~~ ^{subset} elements of a big universe $S \subseteq U = \{0, \dots, u-1\}$ in a small table $T = \{0, \dots, t-1\}$.

We use randomness to avoid bad worst-case behavior in expectation

Def A class of hash functions \mathcal{H} is

C-universal iff:

$$\forall x, y (x \neq y) \quad \left| \{h \in \mathcal{H} \mid h(x) = h(y)\} \right| \leq C \cdot \frac{|\mathcal{H}|}{u}$$

C-universal classes imply: if $h \in \mathcal{H}$ is chosen randomly:

$$\forall x, y, (x \neq y) \quad \mathbb{P}[h(x) = h(y)] \leq \frac{C}{u}$$

(for h chosen uniformly at random \leftarrow $\frac{1}{|\mathcal{H}|}$)

We pick m such that $\frac{c}{m} \in O(1)$

and study the collisions with $x \in U$
 (Chain-length)

\downarrow collisions of x

$$Z_x = |\{h(x) = h(y) \mid y \in U\}| = \sum_{y \in U} \begin{cases} 1 & \text{if } h(x) = h(y) \\ 0 & \text{else} \end{cases}$$

$$\begin{aligned} E(Z_x) &= \sum_{y \in U} \delta_{h(x), h(y)} = \sum_{y \in U} \Pr[h(x) = h(y)] \\ &= 1 + \sum_{\substack{y \in U \\ y \neq x}} \Pr[h(x) = h(y)] \end{aligned}$$

$$\leq 1 + \frac{c}{m} \cdot n$$

$$= O(1)$$

□

Example for a 2-universal ~~hash~~ family of Hash-functions:

$$h_{a,b}(x) = ((ax + b) \bmod p) \bmod m$$

with p prime

Proof: Textbook (e.g. Cormen)

$W_{arb}(x)$ uses: compact representation $O(1)$
: $O(1)$ evaluation

↳ good hash-function

Summary: Universal hashing with chaining
needs $O(1)$ time / $O(p)$ in
expectation / amortized

drawback: worst case for a "bad" sequence
we may have long query/insert
times \uparrow
 $\Theta\left(\frac{\log n}{\log \log n}\right) =$ longest chain
 \downarrow

Solutions for the static case

~~[Mitzenmacher '96]~~ [Azar et al, '94] [Mitzenmacher '96]

use two (d) hash functions and add the key
in the currently shorter list.

↳ longest chain $\Theta(\log \log n)$

$$\left(\Theta\left(\frac{\log \log n}{\log d}\right) \right)$$

Perfect Hashing (Static)

[Fredman, Komlós, Szemerédi '84]

For a fixed set of n keys find a hash-function h that produces no collisions.

As usual:

- h has compact representation
- h can be evaluated in $O(1)$

$z(h) = \#$ of all collisions of h

$$= \sum_{\substack{x, y \in S \\ x \neq y}} \delta_{h(x), h(y)}$$

Assume h is \mathcal{R} -universal

$$\begin{aligned} \mathbb{E}(z(h)) &= \sum_{\substack{x, y \in S \\ x < y}} \mathbb{E}(\delta_{h(x), h(y)}) = \sum_{\substack{x, y \in S \\ x < y}} \mathbb{P}[h(x) = h(y)] \\ &\leq \binom{n}{2} \frac{2}{m} \leq \frac{n^2}{m} \end{aligned}$$

$\leq \frac{c}{m}$
 $\frac{n^2}{m} \leq \frac{c}{m} \iff n^2 \leq c$ (*)

\Rightarrow For table of size $2n^2$: $\mathbb{E}(z(h)) \leq \frac{1}{2}$

Aside: Markov inequality:

X non-neg Random variable:

$$\mathbb{P}[X > a \cdot \mathbb{E}(X)] \leq \frac{1}{a}$$

Thus: $\mathbb{P}[z(h) \geq 2 \mathbb{E}(z(h))] \leq \frac{1}{2}$

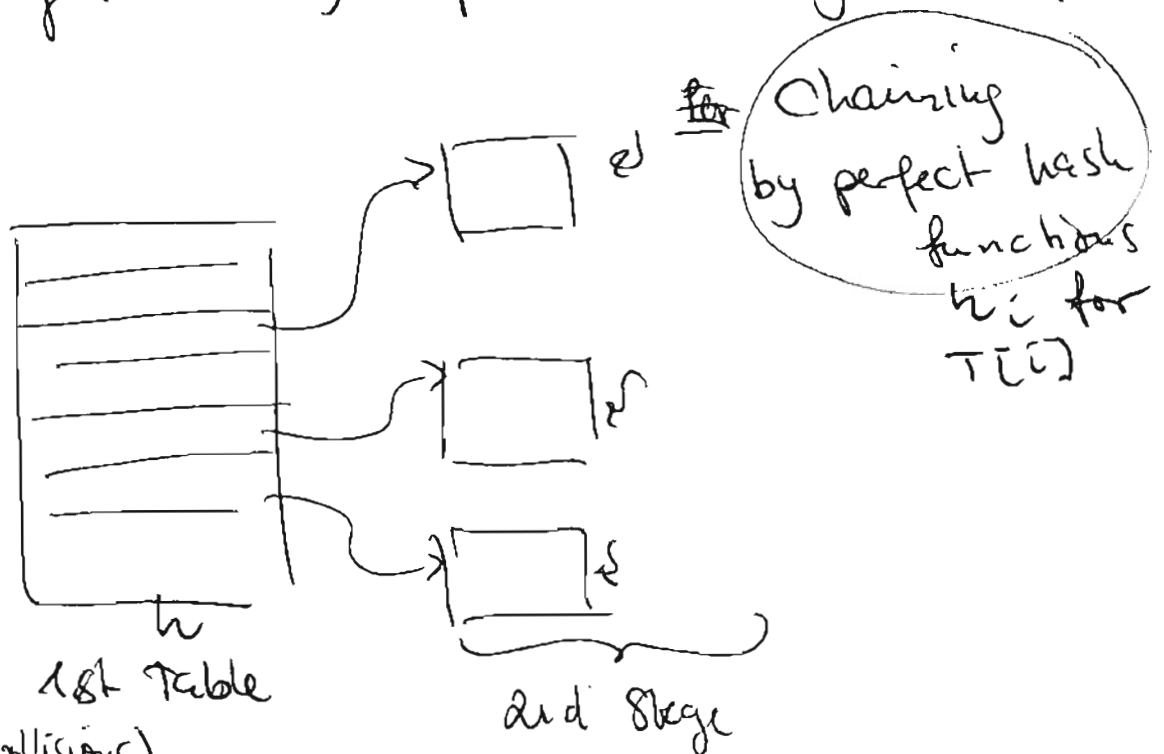
$$\Rightarrow \mathbb{P}[z(h) > 1] \leq \frac{1}{2}$$

\hookrightarrow We can find a "perfect" hash-function by picking h randomly and check if h is perfect. If not, pick again.

Table size: $O(n^2)$

We can get $O(n)$ space using 2-stage

hashing:



$\Rightarrow S$ is partitioned into $S_0, S_1, S_2, \dots, S_{s-1}$
 with: $S_i = \{x \in S \mid h(x) = i\}$

$$Z(u) = \sum_{i=0}^{s-1} |\{ (x, y) \mid x, y \in S, x < y, h(x) = h(y) = i \}|$$

$$= \sum_{i=0}^{s-1} \binom{|S_i|}{2}$$

Hence $E\left(\sum_{i=0}^{s-1} |S_i|^2\right) = E\left(\sum_{i=0}^u 4|S_i|\right) + 2E\left(\sum_{i=0}^u \binom{|S_i|}{2}\right)$

$$2\binom{u}{2} = \frac{u(u-1)}{2} = \frac{1}{2}u^2 - u$$

$$= u + 2E\left(\sum_{i=0}^{s-1} \binom{|S_i|}{2}\right)$$

$$= u + 2 \frac{u^2}{2m} \quad (\text{Due to } (*))$$

$$= 2u \quad \text{for } m = 2u$$

$\Rightarrow O(u)$ space

Due to Markov's inequality:

$$P\left(\sum |S_i|^2 \geq 4u\right) \leq \frac{1}{2}$$

\Rightarrow (1) Choose h until $\sum |S_i|^2 < 4u$

\Rightarrow (2) Choose $\{h_i\}$ until $\sum h_i$ is perfect \square

Cuckoo-Hashing

[Pagh-Rodler 2004]

Def A class of hash-functions \mathcal{H} is said to be (c, k) -universal, iff

$\exists x_1, x_2, \dots, x_k \in \mathcal{U}$ (distinct)

$\exists a_1, a_2, \dots, a_k \in \mathcal{T}$ (distinct):

$$\Pr [h(x_1) = a_1 \wedge h(x_2) = a_2 \wedge \dots \wedge h(x_k) = a_k] = \frac{c \cdot |\mathcal{U}|}{|\mathcal{T}|}$$

For cuckoo hashing we need $2 \cdot (c \cdot \log u)$ universal hash functions, f, g

- f, g can be found ~~in~~ (representable ^{ed} in $O(\log u)$ space and evaluated in $O(1)$)

[Biegel '89]

- Idea Cuckoo hashing: insert in $x \rightarrow \mathcal{T}[f(x)]$

, if there is some $y \in \mathcal{T}[f(x)]$, insert

y in $\mathcal{T}[g(y)]$ or $\mathcal{T}[f(y)]$

- 7- • Stop after $O(\log u)$ evictions

Find (x) : Check $T[f(x)]$ & $T[g(x)]$

Delete (x) : If $x \in T[f(x)]$ or $T[g(x)]$
remove x from the table

insert (x) : $i \leftarrow 1$

while ($i \leq \text{maxloop}$)

if $T[f(x)] = \text{empty}$

then $T[f(x)] \leftarrow x$

else

return
swap $x \leftrightarrow T[f(x)]$

if $T[g(x)] = \text{empty}$

then $T[g(x)] \leftarrow x$

return

else swap $x \leftrightarrow T[g(x)]$

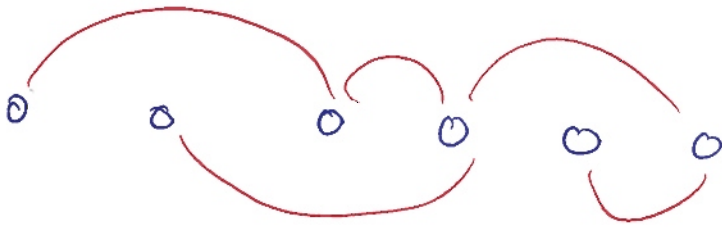
$i \leftarrow i + 1$

Rehash with new f/g
insert (x)

$$\text{maxloop} = O(\log u) = 6 \log u$$

What is the running time of $\text{INSERT}(x)$?

Visualization: cuckoo-graph



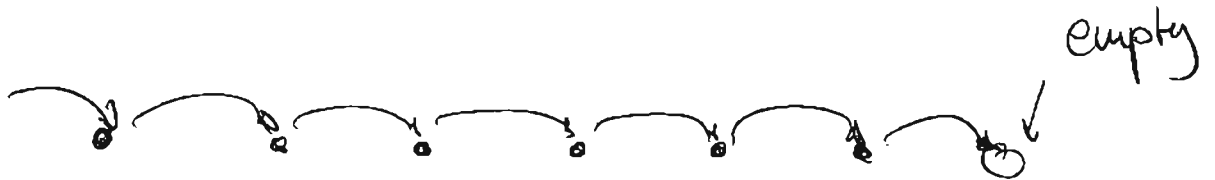
nodes =
cells $\in T$

$$\text{edge } (i, j) \Leftrightarrow \exists x \in \mathbb{R} : f(x) = i, g(x) = j$$

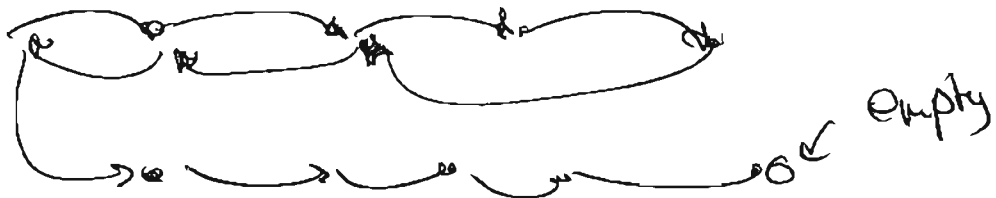
Let n be the # of elements in the table
Assume that $n = 4m$

We have 3 types of situations when inserting

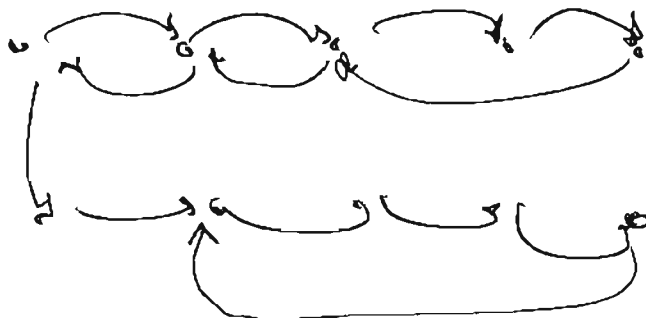
Case A:
(Path)



Case B:
(1 cycle)



Case C:
(2 cycles)



no empty
cell

Case A: $(c \leq 1)$

$\mathbb{P}[T[\varphi(x)] \text{ is occupied}]$

$$= \mathbb{P}(\exists y \in S \setminus \{x\} : \varphi(x) = \varphi(y) \text{ or } \varphi(x) = g(y))$$

$$= \sum_{\substack{y \in S \\ y \neq x}} (\mathbb{P}[\varphi(x) = \varphi(y)] + \mathbb{P}[\varphi(x) = g(y)]) \leq 2 \cdot \frac{uc}{m} = \frac{2uc}{m} = \frac{c}{2}$$

$\mathbb{P}[2 \text{ evictions}] =$

$$\sum_{y_1, y_2} \mathbb{P}[\varphi(x) = \varphi(y_1) \wedge g(y_1) = \varphi(y_2)]$$

$$+ \mathbb{P}[\varphi(x) = g(y_1) \wedge \varphi(y_1) = \varphi(y_2)]$$

$$+ \mathbb{P}[\varphi(x) = \varphi(y_1) \wedge g(y_1) = g(y_2)]$$

$$+ \mathbb{P}[\varphi(x) = g(y_1) \wedge \varphi(y_1) = g(y_2)] \leq 4 \cdot \frac{c}{m^2} = \frac{c}{4}$$

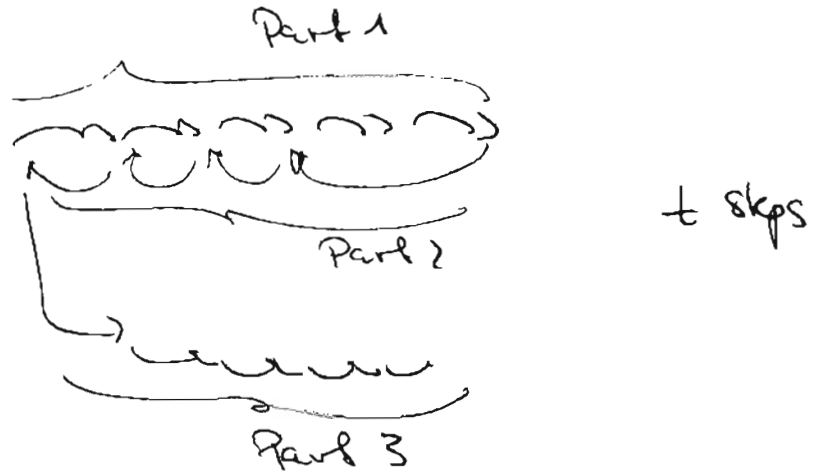
$$\mathbb{P}[t \text{ evictions}] \leq \frac{c}{2^t}$$

$$\Rightarrow \mathbb{E}(\text{evictions in case A}) = \sum_{t=0}^{\infty} (t+1) \frac{c}{2^t} = O(1)$$

$$\mathbb{P}(\text{release in case A}) = 2^{-\log_u n} \leq 2^{-\log_u^2} = O\left(\frac{1}{u^2}\right)$$

Case B:

Observation



Because $\text{Part 2} < \text{Part 1}$: $\text{Part 1 or Part 3} \geq \frac{t}{3}$

$\Rightarrow \mathbb{P}[\text{1 loop with } t \text{ steps}] <$

$$\mathbb{P}[\text{path of length } \frac{t}{3} \text{ from } x] = 2^{-\frac{t}{3}}$$

$\Rightarrow \mathbb{P}(\text{runny tree case B}) \leq \sum (t+1) 2^{-\frac{t}{3}} = O(1)$

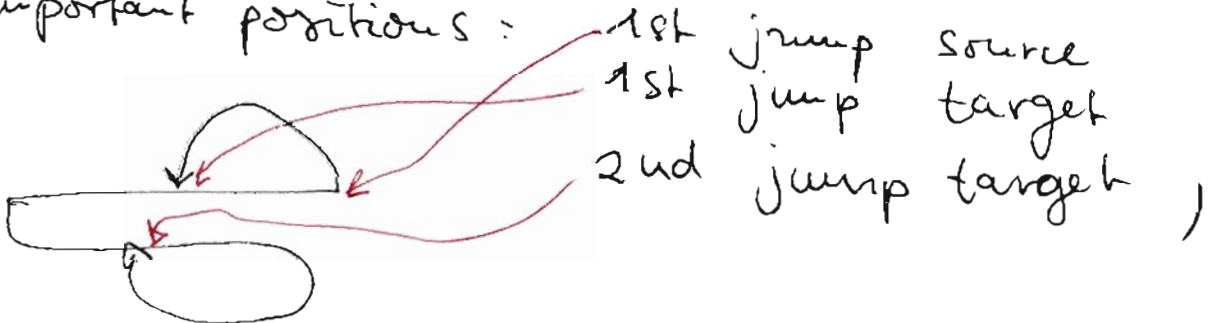
$$\mathbb{P}(\text{rehash}) = 2^{-\frac{6}{3} \log u} = O\left(\frac{1}{4^2}\right)$$

Case C (via counting)

How many 2-loop scenarios length t distinct edges nodes

: combinatorial 2-loop graphs: $\leq t^3$

(3 important positions:



How many ways to place a loop in the table? : $m^{\cancel{t}-1} = (4u)^{t-1}$ (position of x is assigned)

How many ways to place keys in the cells? : u^{t-1}

\Rightarrow in total $t^3 u^{t-1} (4u)^{t-1}$ configurations

How many circuits graphs are there?
with t edges

\forall edge $(u, v) \exists f, g : f(y) = u \quad g(y) = v$
or $g(y) = u \quad f(y) = v$

$\Rightarrow \forall$ edge : $\frac{2}{2} = \frac{(4u)^2}{2}$ choices

\Rightarrow in total $\frac{(4u)^{2t}}{2^t}$

$$\begin{aligned} \Rightarrow P[\text{we have case C}] &= \frac{2^t}{(4u)^{2t}} \cdot t^3 u^{t-1} (4u)^{t-1} \\ &= \frac{1}{8u^2} \frac{t^3}{2^{t-1}} \quad O(1) \end{aligned}$$

$$\begin{aligned} \Rightarrow P[\text{we have case C}] &\leq \sum_{t=2}^{\infty} \frac{t^3}{8u^2 2^{t-1}} = \frac{1}{8u^2} \sum_{t=2}^{\infty} \frac{t^3}{2^{t-1}} \\ &= O\left(\frac{1}{u^2}\right) \end{aligned}$$

$$\Rightarrow P[\text{rehash}] = O\left(\frac{1}{u^2}\right) \begin{cases} \text{long case A, B} \\ \text{case C} \end{cases}$$

$$\Rightarrow P[\text{rehash after } u \text{ insertions}] = O\left(\frac{1}{u}\right)$$

$$\Rightarrow P[\text{rehash is successful}] = 1 - O\left(\frac{1}{u}\right)$$

↳ successful after $O(1)$ tries

~~\Rightarrow Rehashing takes $O(u)$ time in expectation~~

~~↳ amortizable over ku~~

$$\Rightarrow \text{Case C: costs } \underbrace{O(\log u)}_{\substack{\uparrow \\ \text{max loop steps}}} + \underbrace{O(u)}_{\substack{\uparrow \\ \text{rehash}}}$$

$$\Rightarrow \text{in expectation: } O\left(\frac{1}{u^2} \cdot O(u)\right) = O\left(\frac{1}{u}\right)$$

$$\Rightarrow \text{Amortized time INSERT: } O(1) + O\left(\frac{1}{u}\right) = O(1)$$

in expectation,
amortized

□