

Perfect Hashing, Cuckoo Hashing



Birthday (non)Paradox

- **Q:** How many people must be in a room for it to be more likely than not that two have ***the same birthday***?
- **A:** 23. \mathcal{E} = event that all birthdays are distinct.
 - $\Pr(\mathcal{E}) = 1 \left(1 - \frac{1}{365}\right) \left(1 - \frac{2}{365}\right) \dots \left(1 - \frac{22}{365}\right) < 1/2.$
- $S \subset [u]$ is a set of n keys.
- $h: [u] \rightarrow [m]$ is a uniformly random hash function.
- **Q:** How big must m be for h to be *injective* on S with constant probability?
- **A:** $m = n^2$ will do.
 - And this is true even if h is only pairwise independent!

$$\prod_{i=1}^{n^2} \left(1 - \frac{1}{36r}\right) \approx \prod_{i=1}^{2n} e^{-\frac{i}{36r}} = e^{\sum_{i=1}^{2n} -\frac{i}{36r}} = e^{-\frac{n(n+1)}{36r}}$$

$X = \#\text{ of collisions}$
 (i,j) $i, j \in S$ $P(i) = P(j)$

$$P(X \geq 1) \leq \frac{E(X)}{1} < \frac{1}{2}$$

$$X = X_1 + X_2 - X_{(2)}$$

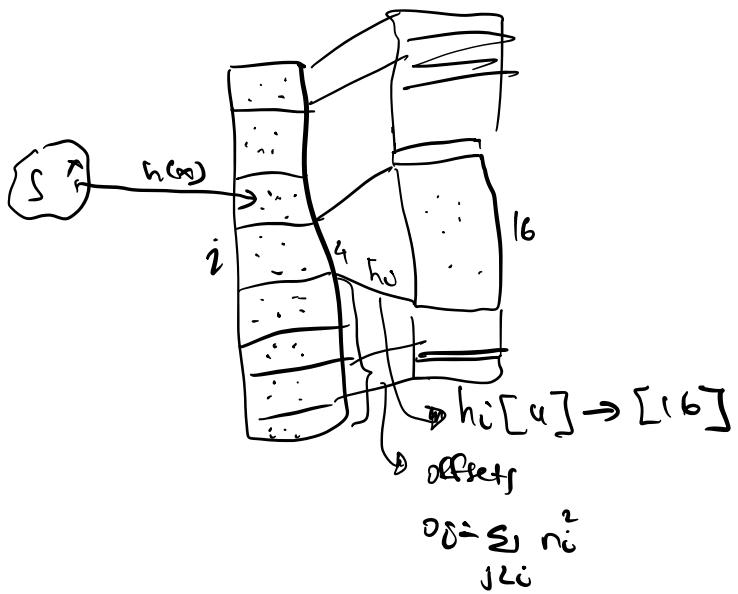
$$E(X) = \sum_{i=1}^{\binom{n^2}{2}} X_i = \binom{n^2}{2} E(X_i) = \binom{n^2}{2} \left(\frac{1}{36}\right) m = n^2$$

$$\xrightarrow{\text{pairwise independence}} = \binom{n^2}{2} P(h(x)=h(y)) \Rightarrow \frac{n(n-1)}{2} < \frac{1}{2}$$

Data Structures

- The **dictionary** problem: maintain a dynamic set S of **keys** (with associated data)
 - Insert(x) : add x to S
 - Delete(x) : delete x from S
 - Lookup(x) : is $x \in S$? (Return ptr. to associated data if yes.)
- **Static case**: S is given in advance; we only need to support Lookup(x) operations.

Perfect hashing in $O(n^2)$ space. 2-wrappers h , alloc $A[0, \dots, n^2 - 1]$
 $S \subseteq [u]$ Store $x \in S$ in $A[h(x)]$. 50% chance h will work, can keep picking h 's. But space is large



back to before, $m=n$

$$X = \# \text{ collisions} \\ (i, j) \quad h(i) = h(j) \quad \text{in } S$$

$$E(X) = \frac{\binom{n}{2}}{n} = \frac{n-1}{2} \quad n_i = |S_i| \quad S_i = \{x \mid h(x) = i\}$$

$$X = \sum_i \binom{n_i}{2} = \sum_{i=1}^n \frac{n_i(n_i-1)}{2} \leq \frac{n}{2} + \sum \frac{n_i^2}{2} \leq n + 2X$$

$$\sum n_i^2 \leq n + 2 \sum \binom{n_i}{2}$$

$$E(\sum n_i^2) \leq (n+2 \sum \binom{n_i}{2}) \leq E(n+2X) \\ \leq n + 2 \left(\frac{n-1}{2} \right) = 2n$$

1) Pick main h
w/ < 50% collisions

2) Then pick each h_i until each i has an injective h_i . ($\approx O(1)$ time)

3) Overall $O(n)$ Space.

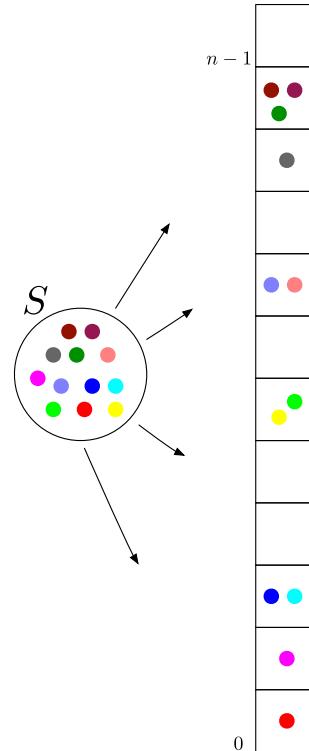
$h: [n] \rightarrow [n]$
 $S_i = \{x \in S \mid h(x) = i\}$
 $n_i = |S_i|$
 $O(n)$

Space: $h, h_1, \dots, h_n, \frac{h}{f_{avg}}$
 $\sum n_i^2$ slots, extra array
 $0, \dots, 0_{n-1}$ offsets

Perfect Hashing in $O(n)$ space

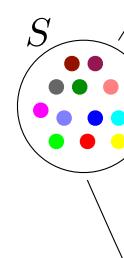
[Ajtai, Komlos, Szemerédi]

- $S = \text{set of } n \text{ keys is } \textit{fixed}.$
- Pick $h: [u] \rightarrow [n]$
 - $E[\# \text{ collisions among } S] = \frac{\binom{n}{2}}{n} = \frac{n-1}{2}$.
- $S_i = \{x \in S \mid h(x) = i\}; \quad n_i = |S_i|.$
- Pick h_0, \dots, h_{n-1} s.t.
 - $h_i : [u] \rightarrow [n_i^2]$
 - and h_i is collision-free on S_i .

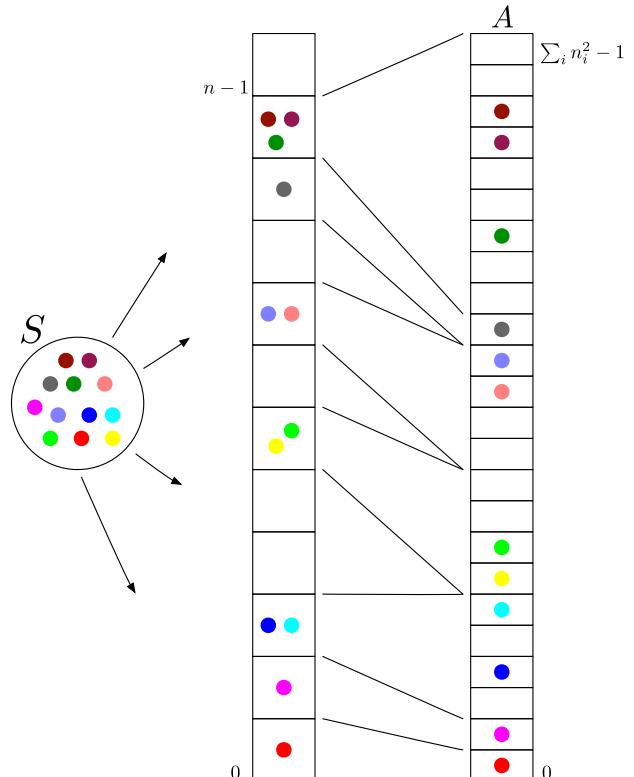


Perfect Hashing in $O(n)$ space

[Ajtai, Komlos, Szemerédi]

- $S = \text{set of } n \text{ keys is fixed}.$
 - Pick $h: [u] \rightarrow [n]$
 - $E[\# \text{ collisions among } S] = \frac{\binom{n}{2}}{n} = \frac{n-1}{2}.$
 - $S_i = \{x \in S \mid h(x) = i\}; \ n_i = |S_i|.$
 - Pick h_0, \dots, h_{n-1} s.t.
 - $h_i : [u] \rightarrow [n_i^2]$
 - and h_i is collision-free on S_i .
 - Allocate array A , $|A| = \sum_i n_i^2$.
 - Store S_i in n_i^2 consecutive cells of A .
 - The data structure consists of h, A ,

h_0, \dots, h_{n-1} , and the offsets $o_i = \sum_{0 \leq j < i} n_j^2$.



Cuckoo Hashing

- Hashing with chaining:
 - $O(m + n)$ space.
 - No need to know maximum value of “ n ” in advance.
 - *Inserts, Deletes, Lookups* in $O(1 + n/m)$ time in expectation.
No worst-case time bound.
- Static perfect hashing:
 - $O(n)$ space.
 - No *Inserts/Deletes* allowed. *Lookups* in $O(1)$ worst-case time.
- **Cuckoo Hashing**
 - $O(n)$ space.
 - *Deletes* and *Lookups* in $O(1)$ worst-case time.
 - *Inserts* in $O(1)$ time *in expectation.*

Cuckoo Hashing

- $n =$ maximum number of elements in set S .
- $r = 2n$
- Allocate arrays $A_1[0..r - 1]$ and $A_2[0..r - 1]$
- Choose **two** hash functions $h_1, h_2 : [u] \rightarrow [r]$

The only invariant:

If $x \in S$ then $A_1[h_1(x)] = x$ or $A_2[h_2(x)] = x$ (but not both)

- $\text{Lookup}(x) : \text{return}(A_1[h_1(x)] == x \text{ || } A_2[h_2(x)] == x)$
- $\text{Delete}(x) : \text{if } (A_1[h_1(x)] == x) \text{ } A_1[h_1(x)] = \text{NULL};$
 $\text{if } (A_2[h_2(x)] == x) \text{ } A_2[h_2(x)] = \text{NULL};$
- $\text{Insert}(x) : \text{the only non-trivial one...}$

Cuckoo Hashing

- $\text{Lookup}(x)$: $\text{return}(\text{A}_1[\text{h}_1(x)] == x \text{ || } \text{A}_2[\text{h}_2(x)] == x)$
- $\text{Delete}(x)$: $\text{if } (\text{A}_1[\text{h}_1(x)] == x) \text{ } \text{A}_1[\text{h}_1(x)] = \text{NULL};$
 $\text{if } (\text{A}_2[\text{h}_2(x)] == x) \text{ } \text{A}_2[\text{h}_2(x)] = \text{NULL};$
- $\text{Insert}(x)$:
 $\text{first check that } x \notin S$
 $\text{z} = x;$
 $\text{while}(\text{should something go here?}) \{$
 $\text{swap}(z, \text{A}_1[\text{h}_1(z)])$;
 $\text{if } z == \text{NULL} \text{ return};$
 $\text{swap}(z, \text{A}_2[\text{h}_2(z)])$;
 $\text{if } z == \text{NULL} \text{ return};$
}

“ $\text{swap}(u, v)$ ” just swaps the values of u and v .

Cuckoo Hashing

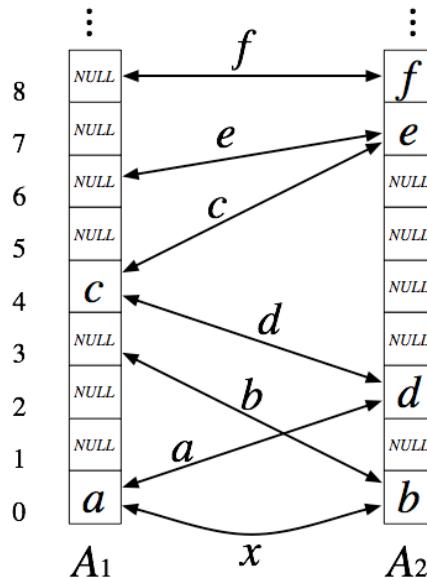
- Suppose a, b, c, d, e, f are already in the table
- $h_1(a) = 0, \quad h_1(c) = 4$
- $h_2(b) = 0, \quad h_2(d) = 2, \quad h_2(e) = 7, \quad h_2(f) = 8$

| | | | | |
|---|---|----------|---|----------|
| | ⋮ | | ⋮ | |
| 8 | | NULL | | <i>f</i> |
| 7 | | NULL | | <i>e</i> |
| 6 | | NULL | | NULL |
| 5 | | NULL | | NULL |
| 4 | | <i>c</i> | | NULL |
| 3 | | NULL | | NULL |
| 2 | | NULL | | <i>d</i> |
| 1 | | NULL | | NULL |
| 0 | | <i>a</i> | | <i>b</i> |
| | | A_1 | | A_2 |

Cuckoo Hashing

- Suppose a, b, c, d, e, f are already in the table
- $h_1(a) = 0, h_1(b) = 3, h_1(c) = 4, h_1(d) = 4, h_1(e) = 6, h_1(f) = 8$
- $h_2(a) = 2, h_2(b) = 0, h_2(c) = 2, h_2(d) = 2, h_2(e) = 7, h_2(f) = 8$
- Suppose we try $\text{Insert}(x)$:

$$h_1(x) = h_2(x) = 0$$



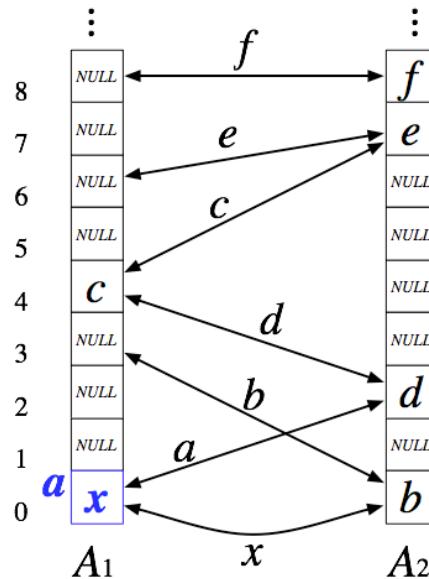
Cuckoo Hashing

- Suppose a, b, c, d, e, f are already in the table
- $h_1(a) = 0, h_1(b) = 3, h_1(c) = 4, h_1(d) = 4, h_1(e) = 6, h_1(f) = 8$
- $h_2(a) = 2, h_2(b) = 0, h_2(c) = 2, h_2(d) = 2, h_2(e) = 7, h_2(f) = 8$
- Suppose we try $\text{Insert}(x)$:

$$h_1(x) = h_2(x) = 0$$

$z = x;$

$\text{swap}(z, A_1[0]);$ (now $z=a$)



Cuckoo Hashing

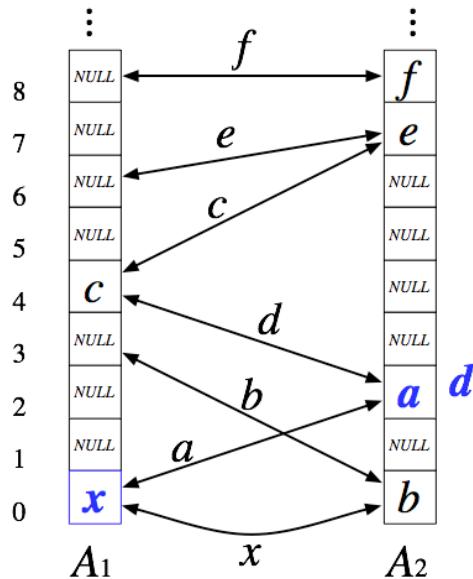
- Suppose a, b, c, d, e, f are already in the table
- $h_1(a) = 0, h_1(b) = 3, h_1(c) = 4, h_1(d) = 4, h_1(e) = 6, h_1(f) = 8$
- $h_2(a) = 2, h_2(b) = 0, h_2(c) = 2, h_2(d) = 2, h_2(e) = 7, h_2(f) = 8$
- Suppose we try $\text{Insert}(x)$:

$$h_1(x) = h_2(x) = 0$$

$z = x;$

$\text{swap}(z, A_1[0]);$ (now $z=a$)

$\text{swap}(z, A_2[2]);$ (now $z=d$)



Cuckoo Hashing

- Suppose a, b, c, d, e, f are already in the table
- $h_1(a) = 0, h_1(b) = 3, h_1(c) = 4, h_1(d) = 4, h_1(e) = 6, h_1(f) = 8$
- $h_2(a) = 2, h_2(b) = 0, h_2(c) = 2, h_2(d) = 2, h_2(e) = 7, h_2(f) = 8$
- Suppose we try $\text{Insert}(x)$:

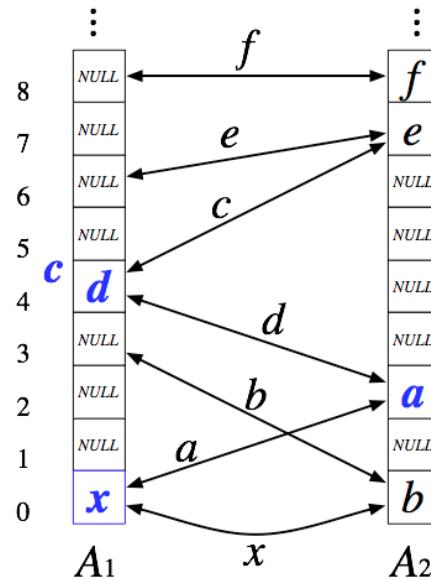
$$h_1(x) = h_2(x) = 0$$

$z = x;$

$\text{swap}(z, A_1[0]);$ (now $z=a$)

$\text{swap}(z, A_2[2]);$ (now $z=d$)

$\text{swap}(z, A_1[4]);$ (now $z=c$)



Cuckoo Hashing

- Suppose a, b, c, d, e, f are already in the table
- $h_1(a) = 0, h_1(b) = 3, h_1(c) = 4, h_1(d) = 4, h_1(e) = 6, h_1(f) = 8$
- $h_2(a) = 2, h_2(b) = 0, h_2(c) = 2, h_2(d) = 2, h_2(e) = 7, h_2(f) = 8$
- Suppose we try $\text{Insert}(x)$:

$$h_1(x) = h_2(x) = 0$$

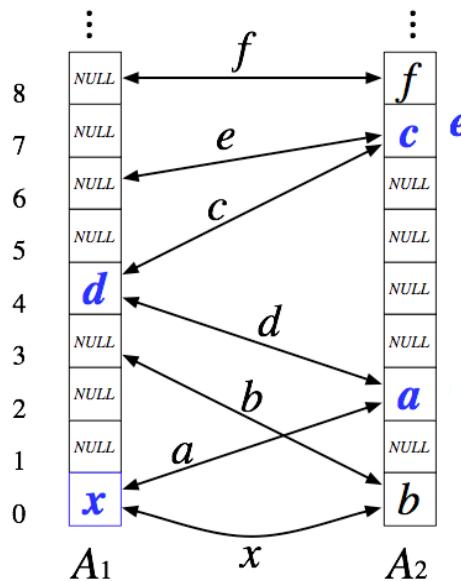
$z = x;$

$\text{swap}(z, A_1[0]);$ (now $z=a$)

$\text{swap}(z, A_2[2]);$ (now $z=d$)

$\text{swap}(z, A_1[4]);$ (now $z=c$)

$\text{swap}(z, A_2[7]);$ (now $z=e$)



Cuckoo Hashing

- Suppose a, b, c, d, e, f are already in the table
- $h_1(a) = 0, h_1(b) = 3, h_1(c) = 4, h_1(d) = 4, h_1(e) = 6, h_1(f) = 8$
- $h_2(a) = 2, h_2(b) = 0, h_2(c) = 2, h_2(d) = 2, h_2(e) = 7, h_2(f) = 8$
- Suppose we try $\text{Insert}(x)$:

$$h_1(x) = h_2(x) = 0$$

$z = x;$

$\text{swap}(z, A_1[0]);$ (now $z=a$)

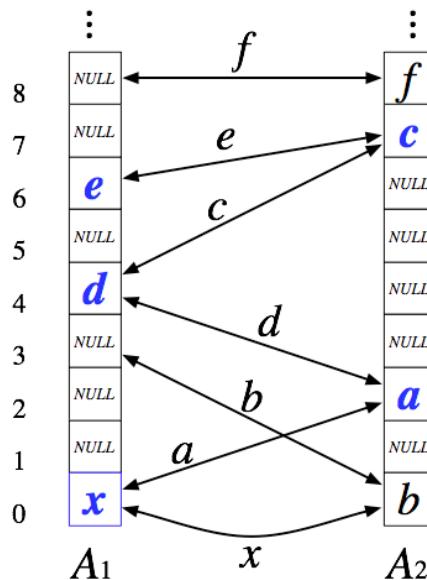
$\text{swap}(z, A_2[2]);$ (now $z=d$)

$\text{swap}(z, A_1[4]);$ (now $z=c$)

$\text{swap}(z, A_2[7]);$ (now $z=e$)

$\text{swap}(z, A_1[6]);$

$z==\text{NULL}$ and we're done





What's the worst that could happen?

- Insert(x) : (first check that $x \notin S$)

```
z = x;  
while(should something go here?) {  
    swap(z, A1[h1(z)]);  
    if z==NULL return;  
    swap(z, A2[h2(z)]);  
    if z==NULL return;  
}
```

What's the worst that could happen?

- $\text{Insert}(x)$: (first check that $x \notin S$)

```
z = x;
```

```
while(we haven't seen a "closed loop") {  
    swap(z, A1[h1(z)]);  
    if z==NULL return;  
    swap(z, A2[h2(z)]);  
    if z==NULL return;  
}
```

pick new hash functions h_1, h_2

rebuild the whole data structure ←

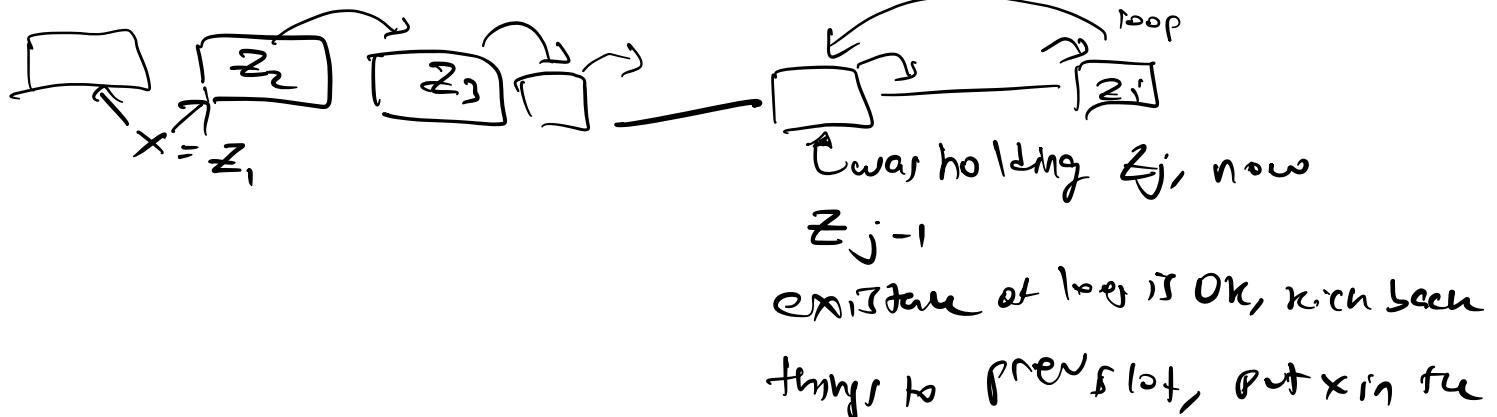
This takes $\Omega(n)$ time! We need to
guarantee this doesn't happen too often

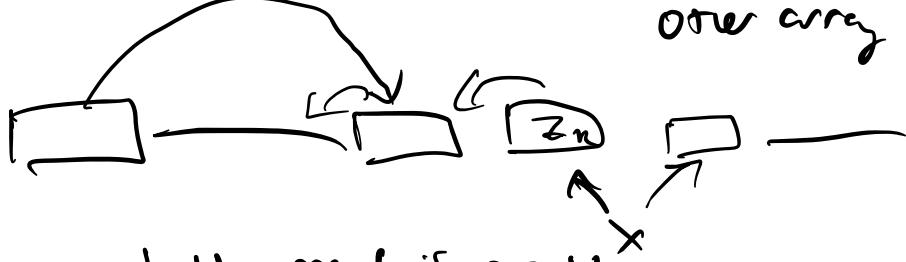
The analysis

- Review of parameters
 - $S \subseteq [u]$ = set of keys in data structure,
 - u = universe size.
 - $|S| \leq n$ at all times.
 - $r = 2n$. A_1 and A_2 both have length r .
 - Assume the ***IDEAL hash model***:

h_1, h_2 are selected uniformly at random

from ALL functions from $[u] \rightarrow [r]$





loop on both arrays is a problem.

"Closed loop" vertices - array keys
edges - keys

\equiv a minimal graph w/ v keys / edges

$v-1$ array keys / vertex



Error Prob

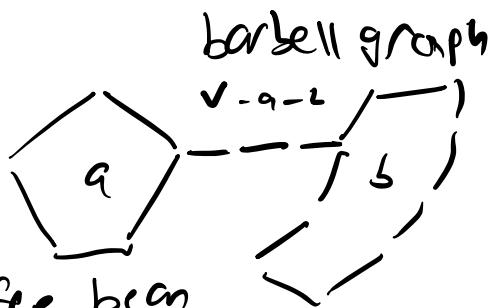
$P(C)$ finding a closed loop involving X)

$$= \sum_{v \geq 3} \sum_{\substack{\text{closed loops} \\ c \text{ w/ size } v}} P(C \text{ exists})$$

1. Choose c ($\leq v^2$) topologically
2. representation in array

Since each vertex is an array cell

3. Choose Keys & Assignment to edges



coffee bean



$$9 + 1 - v$$

$\leq r^{v-1}$ choices of cells each time
 $\binom{r}{v-1} v! \cdot r^{-2v}$ assign edges
 choose keys

$= \sum_{v \geq 3} v^2 \binom{r}{v-1} r^{v-1} v! \cdot r^{-2v}$

$\theta(h_1, h_2) =$
 $h_1(x) =$
 $h_2(x) =$
 $h_1(e) =$
 $h_2(e) =$

$= r^{-2v}$

$= \sum_{v \geq 3} v^3 \frac{r^{v-1}}{(v-1)!} r^{v-1} v! \cdot r^{-2v}$

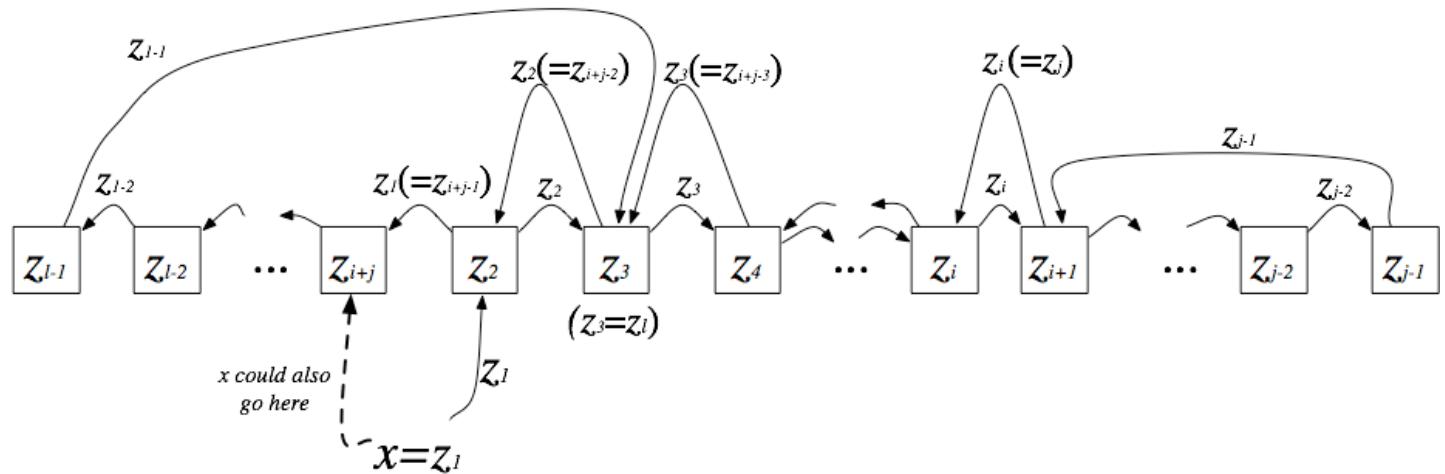
$= \frac{1}{n} \cdot \sum_{v \geq 3} v^3 \left(\frac{r}{n}\right)^v \approx O\left(\frac{1}{n}\right) \underbrace{\left(\sum_{v \geq 3} v^3 \left(\frac{1}{2}\right)^v\right)}_{= O\left(\frac{1}{n^2}\right)}$ since $r=2n$

geometric sum
(converges to constant, $O(1)$)

Upper bound on prob that when you insert an edge, it is the last piece of a closed loop.

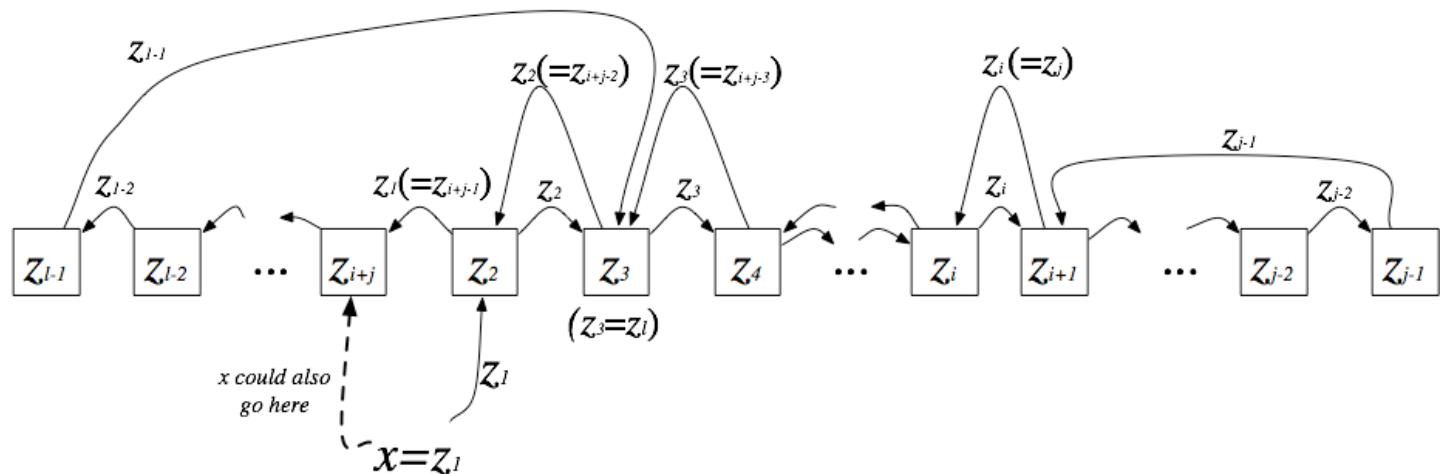
What does a “closed loop” look like?

- Let z_g be the g^{th} value of var. z during an insert.
- Boxes = table cells (alternating A_1 and A_2)
- Boxes labeled with the z_g *originally* there (before the insert)
- Arrow labels = the z_g being relocated during an insert
- By def: $z_j = z_i$ is the *first* key seen twice
- By def: z_l is the *next* key seen twice where $l \geq i+j$.
- If l exists then we have a ***closed loop***.



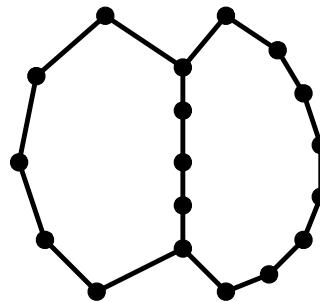
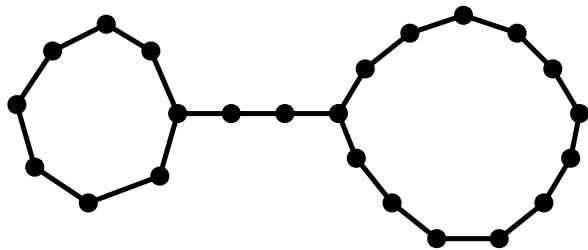
and why is this a bad situation??

- There are $|I|-i-1$ distinct keys ($z_1, \dots, z_{j-1}, z_{i+j}, \dots, z_{|I|-1}$).
- They cannot be put in only $|I|-i-2$ table locations.
- In the analysis “ v ” = $|I|-i-1$ is the number of keys involved in the closed loop.



Closed Loops as Graphs

- Vertices : array cells; Edges : elements in S .
- Closed Loops are minimal graphs with more edges than vertices. They come in two varieties:
- “Bar Bells” and “Coffee Beans”



- How many topologically distinct closed loops are there with v edges and $v - 1$ vertices?

Closed loops are possible, but likely?

- The strategy:
 - (1) Enumerate EVERY POSSIBLE closed loop involving x .
 - (2) Calculate the probability that this closed loop occurs.
 - (3) Sum all probabilities from (2). This is an **upper bound** on the probability of there being ***any*** closed loop involving x .

$$\Pr(\exists \text{ a closed loop}) \leq \sum_{v \geq 3} \sum_{\substack{\text{every closed} \\ \text{loop } C \text{ with } v \text{ keys}}} \Pr(C \text{ exists})$$

What to do now:

Count the number of closed loops with v keys

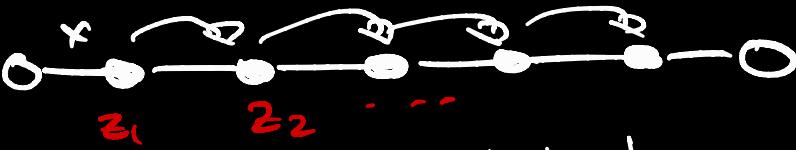
Calculate the probability that such a loop exists.

prob comp

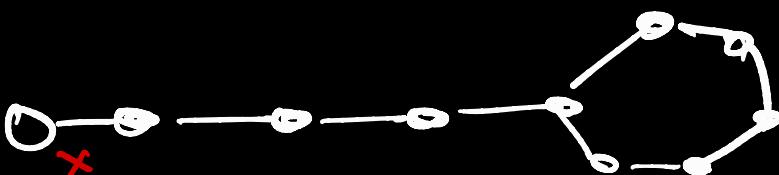
$$\Pr(\exists \text{ a closed loop}) \leq \sum_{v \geq 3} \underbrace{\left(v^2 \binom{n}{v-1} r^{v-1} v! \right)}_{\text{counting component}} (r^{-2v})$$
$$\leq \sum_{v \geq 3} \left(v^2 \frac{n^{v-1}}{(v-1)!} \cdot r^{v-1} v! \right) (r^{-2v})$$
$$\leq O(1) \cdot \sum_{v \geq 3} v^3 n^{v-1} r^{v-1} r^{-2v}$$
$$= O\left(\frac{1}{nr}\right) \cdot \sum_{v \geq 3} v^3 \left(\frac{n}{r}\right)^v = O\left(\frac{1}{n^2}\right) \sum_{v \geq 3} v^3 2^{-v} = O\left(\frac{1}{n^2}\right)$$

Expected time for insert

→ path with $r \geq 1$ (edges) and v vertices



→ lollipop with $r \geq 1$ keys edges & v vertices



→ with \checkmark keys / cells how many of them
shape?

↳ paths → 1

↳ lollipops → r

Expected Time for Insert

$\Pr(\exists \text{ insertion path with size } v \text{ containing } x)$

$$\leq \binom{n}{v-1} r^v (v-1)! \cdot r^{-(2v-1)}$$

$$\leq \left(\frac{n}{r}\right)^{v-1} = \left(\frac{1}{2}\right)^{v-1}$$

$\Pr(\exists \text{ lollipop with size } v \text{ containing } x)$

$$\leq v \binom{n}{v-1} r^v v! \cdot r^{-2v}$$

$$\leq v^2 n^{v-1} r^{-v} = O\left(\frac{v^2}{n}\right) \left(\frac{1}{2}\right)^{v-1}$$

quite small $r \ll n$

→ choose keys $\binom{n}{r-1}$

→ choose cells $\leq r^v$

→ assign keys $\rightarrow (r-1)!$

→ prob + get this
(x second free)

$$\approx r^{-(2v-1)}$$

extra v come from here

Exp search tree

Expected Time for Insert

- $E(\text{insert time} \mid \text{no closed loop})$

$$\begin{aligned}
 &= \sum_{v \geq 1} O(v) \cdot \Pr(\exists \text{ size } v \text{ path/lollipop}) \\
 &= \sum_{v \geq 1} O(v) \cdot \left(\frac{1}{2}\right)^{v-1} \left(1 + O\left(\frac{v^2}{n}\right)\right) \\
 &= O(1)
 \end{aligned}$$

$$\Pr(\#\text{insert swaps} > 2\log n) \leq O\left(\frac{\log^2 n}{n^2}\right).$$

This is when you
 shade the
 search
 for inserting

- Insert(x) :


```
(first check that x ∉ S)
z = x;
while(#swaps ≤ 2log n) {
  swap(z, A1[h1(z)]);
  if z==NULL return;
  swap(z, A2[h2(z)]);
  if z==NULL return;
}
pick new hash functions h1, h2
rebuild the whole data structure
```

The Punchline

- If insertion time is capped at $O(\log n)$, then all events in the analysis involve $O(\log n)$ keys.
 - [Pagh, Rodler 2001] If h_1, h_2 are $O(\log n)$ -wise independent, then insertions take $O(1)$ time in expectation and fail with probability $\tilde{O}(n^{-2})$.
 - [Cohen, Kane 2009] 5-wise independent hash functions are **not enough**. Closed loops appear with high probability.
 - [Open Problem] Is any k -wise independence enough, for some $k = O(1)$?