# Open Addressing



every slot contains max. 1 item

Generalize hash function $h(x)$
to probe sequence $h(x,0), h(x,1), \ldots, h(x,m-1)$

primary location (alt. location)

permutation of $[m]$

## Hashing with Linear Probing

$$h(x,i) := (h(x) + i) \bmod m$$



run

Good news: ① it's simple
② cache-friendly

Bad news: ① SLOW once long runs start forming

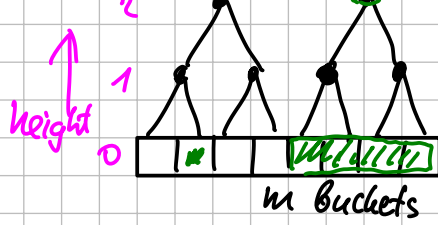More good news: ③ this can be kept under control ☺

Claim: Suppose that $m \geq (1+\epsilon) \cdot n$.
Then $\mathbb{E}[\text{\# probes}]$ is:

① $\Theta(1/\epsilon^2)$ if $h$ is completely random
② $O(1/\epsilon^{13/6})$ for $h$ chosen from 5-indep. family
⑥ $O(1/\epsilon^2)$ for tabulation hashing

③ $\Omega(\log n)$ for some 4-indep. family
④ $\Omega(\sqrt{m})$ for some 2-indep. family
⑤ $\Omega(\log n)$ for multiply-shift

Theorem: Let $m$ be a power of two,
$n \leq m/3$
$h$ be a completely random hash function
$x$ be the item we search for.
Then $\mathbb{E}[\text{\# probes}] \in O(1)$.

Proof: WLOG $n = \frac{1}{3} m$ ± rounding error
(much)
Call the items in the table $x_1 - x_4$.

---

Insert(x): $i \leftarrow 0$
While $B[h(x,i)] \neq \emptyset$: ← or tombstone
$\quad i \leftarrow i+1$
$B[h(x,i)] \leftarrow x$

☺ This succeeds iff $B$ is not full.
☹ Could be very slow.

Find(x): $i \leftarrow 0$
Loop:
$\quad j \leftarrow h(x,i)$
$\quad$ if $B[j] = x$: return TRUE
$\quad$ if $B[j] = \emptyset$: return FALSE
$\quad i \leftarrow i+1$
$\quad$ if $i \geq m$: return FALSE

Delete(x): problematic
replace item by tombstone
after some time rehash all items

$$\text{load} := \frac{\text{\# full buckets}}{\text{\# all buckets}} \in [0,1]$$

**Df:**



$2^t$

block ≡ interval of buckets
below an internal node
of height $t$

hashed vs. stored

a block is critical ≡ # items $\underbrace{\text{hashed}}$ there $> \frac{2}{3} \cdot 2^t$

of size $2^t$ ──── stored in the structure

height 2 1 0

m buckets

Tool: **Chernoff bound for the right tail:**

Let $X_1 - X_k$ be independent random variables with range $\{0,1\}$.

$$X := \sum_i X_i$$
$$\mu := \mathbb{E}[X]$$
$$c > 1$$

Then $\Pr[X > c \cdot \mu] \leq \left(\frac{e^{c-1}}{c^c}\right)^{\mu}$

$e \doteq 2.71828...$

**Lemma:** Let $B$ be a block of size $2^t$.
Then $\Pr[B \text{ is critical}] \leq \left(\frac{e}{4}\right)^{2^t/3} = q^{2^t}$, where $q = \left(\frac{e}{4}\right)^{1/3} < 1$

$\underbrace{\left(\frac{e}{4}\right)}_{<1}$

**Proof:** Indicator random variables:

$$X_i := \begin{cases} 0 \\ 1 \text{ if } h(x_i) \in B. \end{cases}$$

# items hashed to $B$ = $X = \sum_i X_i$

Means: $\mathbb{E}[X_i] = 0 \cdot \Pr[X_i = 0] + 1 \cdot \Pr[X_i = 1]$
$= \Pr[h(x_i) \in B]$
$= \frac{2^t}{m}$ ← independent events

$\mu = \mathbb{E}[X] = \sum_i \mathbb{E}[X_i] = n \cdot \frac{2^t}{m}$

We have $n = \frac{1}{3} m$, so
$\mu = \frac{1}{3} 2^t$.

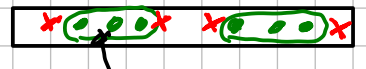$\Pr[B \text{ is critical}]$
$= \Pr\left[X > \frac{2}{3} \cdot 2^t\right]$
$= \Pr[X > 2\mu]$   use Chernoff with $c = 2$
$< \left(\frac{e^1}{2^2}\right)^{\mu} = \left(\frac{e}{4}\right)^{2^t/3}.$ ✓

---

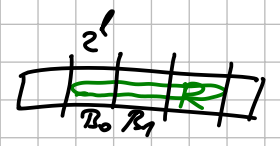**Df:** Run: maximal consecutive set of full buckets
  ☺ the run is preceded and followed by empty bucket
  ☺ an item is hashed to a run ⟺ it's stored in the run



item

**Lemma:** Let $R$ be a run of size at least $2^{\ell+2}$,
  $B_0 - B_3$ be the first 4 blocks of size $2^\ell$ intersected by $R$.
  Then at least one $B_i$ is critical.

$2^\ell$

$B_0 \; B_1 \quad R$

**Proof:** $R$ intersects at least 4 blocks.
$|B_0 \cap R| \geq 1$
$|B_1 \cap R| = 2^\ell$
  ↳ the same
    for $B_2, B_3$

$M := R \cap (B_0 \cup B_1 \cup B_2 \cup B_3)$
$|M| \geq 3 \cdot 2^\ell + 1$

$M$

$B_0 \; B_1 \; B_2 \; B_3 \; \cdots$

what is
stored $M$,
was hashed there

If no $B_i$ is critical: $|M| \leq$ # items hashed to $B_0 - B_3 \leq \frac{2}{3} 2^\ell \cdot 4 = \frac{8}{3} 2^\ell < 3 \cdot 2^\ell$ ⚡

**Lemma:** Let $R$ be the run containing $h(x)$
and $|R| \in [2^{\ell+2}, 2^{\ell+3})$. ⎯ of size $2^\ell$
Then at least 1 of these $R$ blocks is critical:

- the block containing $h(x)$
- 8 blocks before
- 3 blocks after



**Proof:** $|R|$ is between $4 \cdot 2^\ell$ and $8 \cdot 2^\ell$ $\Rightarrow$ $R$ intersects at most 9 blocks
$\rightarrow$ start of $R$ is at most 8 blocks before $h(x)$
& apply the previous lemma.

**Corollary:** Let $R$ be the run containing $h(x)$.
Then $\Pr\left[ |R| \in [2^{\ell+2}, 2^{\ell+3}) \right] \leq 12 \cdot q^{2^\ell}$.

↳ using $\Pr[A \cup B] \leq \Pr[A] + \Pr[B]$

**Finale:**

$$\mathbb{E}[|R|] \overset{\text{def.}}{=} \sum_k k \cdot \Pr[|R| = k] = \left( \sum_{k \leq 3} k \cdot \Pr[\cdots] \right) + \sum_{\ell \geq 0} \sum_{\substack{k \in [2^{\ell+2}, 2^{\ell+3}) }} \overset{\leq 2^{\ell+3}}{(k)} \Pr[|R| = k]$$

↑ run containing $h(x)$

$\in O(1)$

$$2^{\ell+3} \cdot \underbrace{\sum_{k \in \text{Interval}} \Pr[|R| = k]}_{\substack{\Pr[|R| \in \text{Interval}] \\ \leq 12 \cdot q^{2^\ell} \\ \text{Coroll.}}}$$

$$\rightarrow \leq \sum_{\ell \geq 0} 2^{\ell+3} \cdot 12 \cdot q^{2^\ell} = 8 \cdot 12 \cdot \sum_{\ell \geq 0} 2^\ell \cdot q^{2^\ell}$$

$$\leq \boxed{\sum_{t \geq 0} t \cdot q^t} \neq \text{const.}$$

converges as infinite sum for every $q \in (0,1)$

So $\mathbb{E}[|R|] \leq$ some constant. ∎