**Recap:** Searching

- Linear search — $O(n)$ search — $O(1)$ inset an ele
- Binary search — $O(\log n)$ search — $O(n)$ inset an ele.
- Hashing

$\underline{n} \to$ no. of strings in arr.

$O(1)$ search
$O(\log n)$ search

$\underline{O(1)}$ inset time. ✓
$O(\log n)$ inset time ✗

arr = [ "sagan", "amit", "jatin", "anish" ]

| | | | |
|---|---|---|---|
| a → 1 | k 11 | t 20 |
| b → 2 | l 12 | u 21 |
| c — 3 | m 13 | v 22 |
| d 4 | n 14 | w 23 |
| e 5 | o 15 | x 24 |
| f 6 | p 16 | y 25 |
| g 7 | q 17 | z 26 |
| h 8 | r 18 | |
| i 9 | s 19 | |
| j 10 | | |

$n = \underline{10}$ elements $\Rightarrow$ capacity of array

proxy index

"Sagan" $\longrightarrow$ $19 + 1 + 7 + 1 + 18 = \underline{46} \% n = 6$

"amit" $\longrightarrow$ $1 + 13 + 9 + 20 = \underline{43} \% n = 3$

↳ **Insertion logic.**

(assume it is possible to insert)

string $\longrightarrow$ proxy index $\longrightarrow$ go place it.

$\underline{O(1)}$ insert (theoretically)

table

| | |
|---|---|
| 0 | |
| 1 | "anish" |
| 2 | |
| 3 | "amit" |
| 4 | |
| 5 | ⋮ |
| 6 | "sagan" |
| 7 | ⋮ |
| 8 | |
| 9 | |

mapping
ASCII
unicode
custom

**Searching**

$\underline{O(1)}$ search

search ("sagan")

$\longrightarrow$ $46 \% 10 \Rightarrow \boxed{\underline{6}}$

if table[proxy_index] == key:
return true

anish $= \dfrac{1 + 14 + 9 + 19 +}{8} = 51 \% 10$
$= \boxed{1}$

**Terminology:** "str" $\longrightarrow$ hash table

hash ⇐ Method you figure out index with

hash fn.

proxy index $\longrightarrow$ hash value

**Logic is great:** $O(1)$ search $O(1)$ insertion time $\Rightarrow$ Perfect idea. ✗

① What if 2 strings have value 96 and 46 ? table size = 10.

COLLISION. $\longrightarrow$ Resolve this.

# Optimization are required at every step :

**①  What should be the size of the table ?**

Probability :

$$P\left(\text{collision} \mid \#ele = 0\right) = 0 \qquad \underline{\text{occassional}}$$

$$P\left(\text{collision} \mid \#ele = 1\right) = \frac{1}{n} \qquad \leftarrow \quad \underline{\text{low}}$$

$$P\left(\text{collision} \mid \#de = 2\right) = \frac{2}{n} \qquad P(\text{collision}) \propto \frac{1}{n}$$
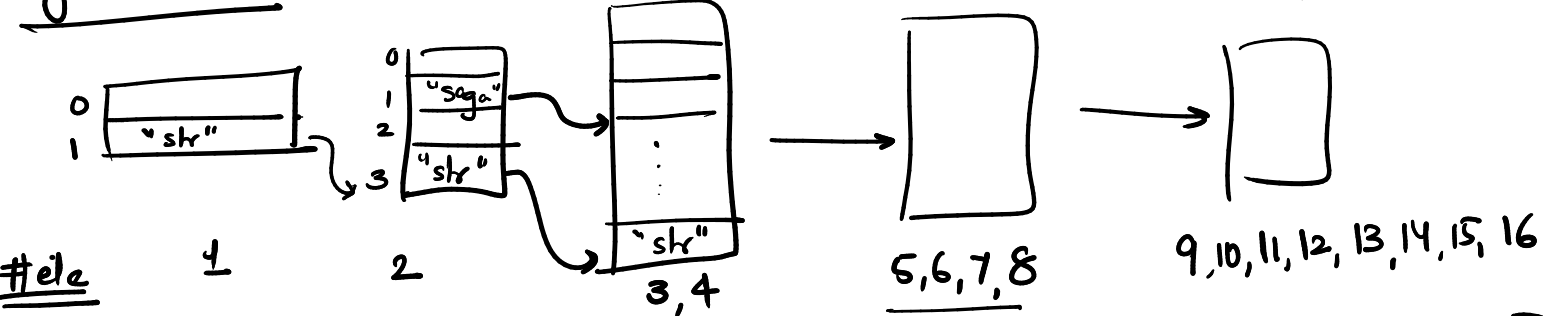
**Q.** I will handle occassional collisions some other way.
After how many elements in array/table of size n, $\rightarrow \frac{n}{2}$
chances of collision are high ?

$$\text{Till } n/2 \text{ elements} \implies P(\text{collision}) < 50\%.$$

most of the times/
atleast 50% of space
in array will be
wasted.

logic on size.

#ele          1          2          3,4          5,6,7,8          9,10,11,12,13,14,15,16

capacity   2 $\rightarrow$ 4 $\rightarrow$ 8 $\rightarrow$ 16 $\rightarrow$ 32 $\rightarrow$ .⑥④

$$\frac{15\% \cdot 2}{①} \qquad \frac{15\% \cdot 4}{\hookrightarrow ③} \qquad 15\% \cdot 8 \Rightarrow ⑦$$

#ele.  ① ② 3 ④ 5 6 7 ⑧ 9 10 11 12 13 14 15
⑯ .... ㉜ $\implies$ When you double the size, #empty
cells increase $\therefore P(\text{collision}) \downarrow$

When you double, you indeed calc all the proxy values of prev
ele again!

T.C Insertion (assuming doubling logic).

Add $\underline{N}$ elements.
$\rightarrow$ no.q elements
$n \rightarrow$ capacity.

All the instances when you recomputed the hash values:

$\underline{\underline{1}}$  $\underline{\underline{2}}$  3  $\underline{\underline{4}}$  5  6  7  $\underline{8}$  .....  N-1  $\underline{N}^{*}$

1  2  ×  4  ×  ×  ×  8                    #work          #had to recompute hashes

Total work
to add      =  $1 + \underline{2} + 4 + ...\underline{\frac{N}{2}} + \boxed{N}$
N elements
in this log.z  = $N \left( \frac{1}{N} + \frac{2}{N} + \frac{4}{N} + .... + 1 \right)$

= $N \left( \underline{\frac{1}{}} + \frac{1}{2} + .... + \frac{4}{N} + \frac{2}{N} + \frac{1}{N} \right)$    converges.

$\leqslant N \left( 1 + \frac{1}{2} + \frac{1}{4} + .... \right)$    Infinite sum GP
$S = \frac{a}{1-r}$

$< N \left( \frac{1}{1-\frac{1}{2}} \right)$

# Recomputates $< \underline{2N}$    to add  N  elements.

To add  N  elements    TC $\Rightarrow O(2N) \sim O(\underline{\underline{N}})$

To add a single element $\longrightarrow$ $\dfrac{O(N)}{N}$ ~ $\underline{O(1)}$     Insertion per element.

Amortization $\longrightarrow$ Vectors in C++     Arraylist in Java     List in Py

Statement : If the indices are proxy $\Rightarrow$ they are bound to change with size of array, I can't return that.

"sagan" $\longrightarrow$ $\underline{\underline{6}}$   X   $\longrightarrow$ Next step I am not so sure of this 6.

bool outcome $\longrightarrow$ present or not.

② $\underline{x \% \,\textcircled{n}}$   there is a uniform distribution.

random "string" $\Longrightarrow$
0
1
2
3
4
⋮

Abstract algebra
$\downarrow$
Group Theory
$\quad\searrow$ Galios theory

$\boxed{n \Rightarrow factors}$

| | | |
|---|---|---|
| 0 | ✓✓✓ | 1/n |
| 1 | ✓ | 1/n |
| 2 | | |
| 2 | ✓✓ | ⋮ |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | 1/n |

$\underline{n = 120}$

n to be prime $\rightarrow$ no factors

factors of n $\rightarrow$ 2, 3, 4, 5, 6, 8, 10, 12, 15, ....

$\underline{Actual}$

#factors $\uparrow$     higher is bias towards $\uparrow$

$\underline{\underline{n}}$ are not exactly double

$\underline{\underline{O(1)}}$ insertion

n: 2   5   11   23

P(collision) $\downarrow$     if $\underline{\underline{n}}$ is prime $\longrightarrow$ n actually blocks the size in RAM.

③ hash f<sup>n</sup>: ⟶ simple      O(1) time to compute hash value.

sum works okay.                    No loops    No recursion.

more randomness than sum           $\underline{bits}$ ⟹ scramble ⟹ $\underline{hash\ value}$
can come from scrambling bits.

                                   $\dfrac{S}{19}$ — —

There will be occassional
collisions.                        | 10110  10001 _ _ . |

Eg.    Anagrams :      Sagar        Ragas
                         ↓            ↓
                                   sum  $\underset{=}{X}$

2 collision resolutions :   ① Chaining
                            ② Probing
                       $\underset{=}{Proxy}$

$S_1$ = "sagar"    ⟶    6                          6 | ("Sagar", "ragas") |
$S_2$ = "ragas"    ⟶    6

Search    "ragas"  ⟶  $\underset{=}{6}$          Among multiple value,
                                                  if anyone matches key ?

                                                      n = 10
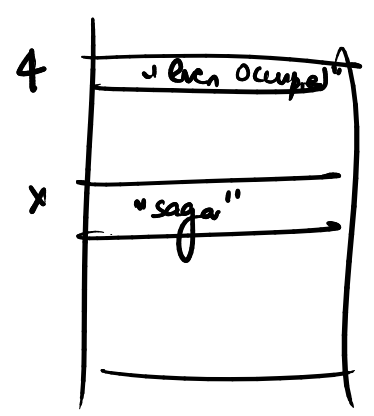
Probing.      str ⟶ h(str) ⟹ x = $\underline{6}$
            ="ragas"                              4 | ✓ even occupied |

probe
f<sup>n</sup>(                                    x | "saga" |
$g(x) = (ax + bx^2) \% m$    a,b ∈ primes

worst
case    $(2(6) + 7(6)^2) \% 10$
practical:  $(12 + 7 \times 36) \% 10$
③         $(12 + 252) \% 10  ⟹  ④.$        $g(x) = a(4) + b(4)^2$

Only drawback
of hashing :
Memory wastage *

Best case scenario
occupy only 50% of table

worst case scenario
occupy only 25% of table

↓

Solve this problem
in future. (Trees)

↓

# Hashing

Search $O(\log n)$     Insertion $O(\log n)$     Extra $= 0$
space

(NO)                                          n items   only n space.
———×   ———×   ——————×   ——————→   ——

Internal :   set.        Java              Py  datatype

C++  STL                 Collections
     set                 HashSet                        set()

                    SET  IS  A  COLLECTION          ↘ uniqueness.
                    OF  UNIQUE OBJECTS.

"saga" ———→ (6)

dup   Insert "saga" ———→ (6) ———→ 6  | "saga" |

(Key) : value        String ———→ int  | count
  |                                     | #times
  ↓
Hash value              FREQUENCY.                    hashtable
 (proxy)                                          → (key , value)

C++:  map          Java: Hashmap        Py : Dict

$$M["sagan"] \longrightarrow M[\underline{\quad}]$$

proxy

$\underline{\quad} \times \underline{\quad} \times \underline{\quad\quad} \times \underline{\quad\quad} \times \underline{\quad\quad}$