| 12 Feb 2025 | Hashing |
|---|---|

(Typeset lecture notes coming soon!)

Balls $\rightsquigarrow$ "keys"

Bins $\rightsquigarrow$ "buckets"

Hashing: assign keys to buckets

<span style="color:red">randomly</span> but <span style="color:red">reproducibly.</span>

## Dictionary   a.k.a. Associative array, Key-value store.

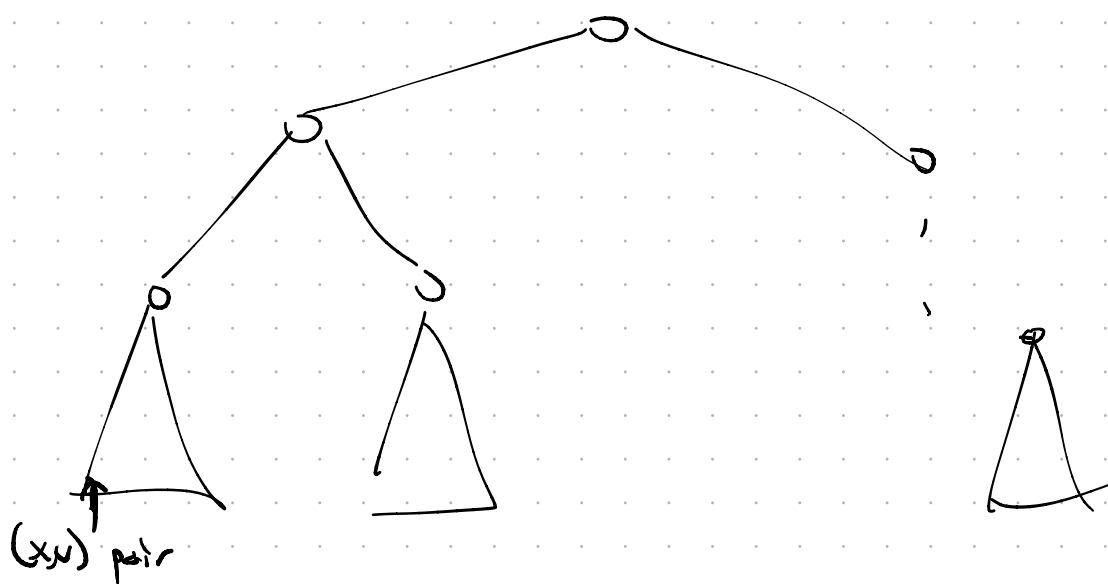Data structure that stores a set of key-value pairs $(x,v) \in X \times V$, with at most one value per key.

Supports:   LOOKUP$(x)$: returns value $V$ stored with $x$
      or   NOT FOUND.

   INSERT$(x,v)$: insert $(x,v)$, overwriting $(x,v')$ if necessary
   DELETE$(x)$: remove $(x,v)$ from data structure
      if any such pair exists.

   a.k.a.   Python dictionaries, Java/C++ HashMaps.

## Deterministic implementation.   Balanced Binary search tree.



$(x,v)$ pair

Tree structured s.t.

a. Every node has equal # Leaves
(within 1) in Left, Right subtrees.

b. All keys in left subtree are
earlier in the ordering of $X$
than the keys in the right subtree.

Assume $m$ key-value pairs organized this way

$\implies$ tree has depth $\lceil \log_2(m) \rceil$

$\implies$ LOOKUP$(x)$ takes $O(\log_2 m)$ time.

Space complexity: $O(m)$ tree nodes

so $O(m)$ space if each key and value
occupies $O(1)$ space.

More generally, the data structure occupies
space $O(S)$ where $S$ is the
min aml. of storage space needed
to store all keys and values.

Hash Table: an array of $n$ buckets.
(Typically $n = O(m)$.)

Each bucket contains a linked list
of key-value pairs. ("chain hashing")

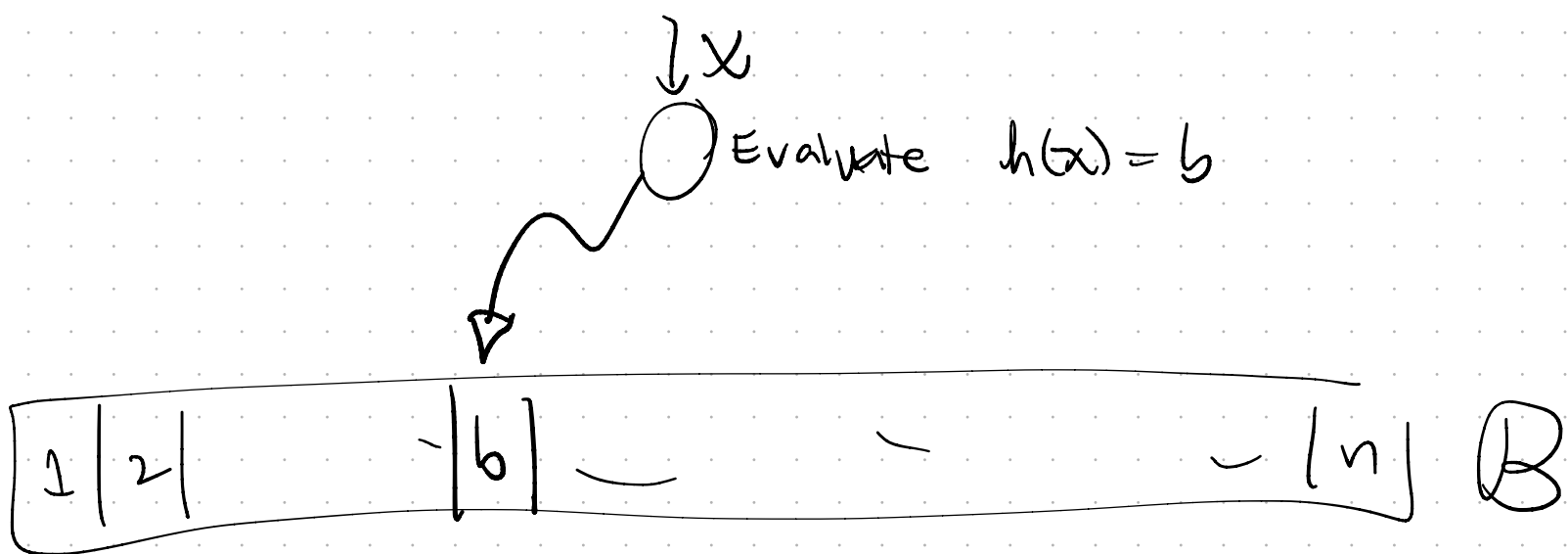A <u>hash</u> <u>function</u> maps keys (randomly

or pseudorandomly) to buckets.

Let $B$ denote {buckets}.

For now assume* (unrealistically) the hash function is
a <u>uniformly random</u> $h: X \to B$,
and that evaluating $h(x)$ takes $O(1)$ time.



$\downarrow x$

Evaluate $h(x) = b$

| 1 | 2 | | ~ | b | ~ | ~ | ~ | n | $B$ |

LOOKUP($x$):     iterates through linked list
                in bucket     $b = h(x)$.
        If $(x, v)$ found,     return $v$.
        Else                   return NOT FOUND.

INSERT($x, v$):   iterates through linked list $b$.
        If $(x, v')$ found, replace with $(x, v)$.
        Else       append $(x, v)$ to linked list

DELETE($x$):   iterates through linked list $b$.
        if $(x, v)$ found     delete it
        else       do nothing.

Actual running time is

$$O\left(1 + \underline{\text{length of linked list stored at } h(x)}\right).$$

random variable depending on $h$
and on the data stored in the table

Linearity of expectation:

$$\mathbb{E}[\text{length of linked list stored at } h(x)]$$

$$= \mathbb{E}\left[\# \; y \in \mathcal{X} \text{ s.t. a pair } (y,w) \text{ is stored in the list at } h(x)\right]$$

$$= \sum_{b \in B} \sum_{y \in \mathcal{X}} \Pr\left[h(x) = h(y) = b \text{ and } (y,w) \text{ is a pair stored in the table}\right]$$

$\leq m$ distinct pairs $(y,w)$ stored in table.

At most one satisfies $y = x$.

For that one $\Pr\left(h(x) = h(y) = b\right)$ is $\frac{1}{n}$.

For the other $\leq m$ pairs $(y,w)$ with $y \neq x$

$$\Pr\left(h(x) = h(y) = b\right) \text{ is } \frac{1}{n^2}.$$

So,

$$\sum_{b \in B} \sum_{y \in \mathcal{X}} \Pr\left[h(x) = h(y) = b \text{ and } (y,w) \text{ is a pair stored in the table}\right]$$

$$\leq \sum_{b \in B} \left[\frac{1}{n} + \frac{m}{n^2}\right]$$

$$= 1 + \boxed{\frac{m}{n}} \longleftarrow \text{"load factor"}$$

Typical implementations use $n \geq m$
$\implies$ load factor $\leq 1$.

Def. A pairwise independent (a.k.a. $2$-universal)
hash family is a probability distribution
on functions $h: X \to B$ satisfying

$$\forall x_1 \neq x_2 \text{ in } X \text{ the pair}$$

$$\left( h(x_1), h(x_2) \right) \text{ is uniformly}$$

$$\text{distributed over } B^2.$$

Example. Suppose $X = B = \{0, 1, 2, \ldots, p-1\}$
and $p$ is prime.
  Sample random $h$ by sampling
coefficients $a, b \in \{0, \ldots, p-1\}$ at random
and defining $h(x) \equiv ax + b \pmod{p}$