

ADVANCED

overall tree
= binary tree of mini trees

mini trees

micro trees

actual nodes

$\frac{1}{7} \lg n$ nodes

$\lg n$ nodes

n nodes



DATA STRUCTURES

5

Integer Data Structures

Advanced Data Structures · Summer 2026

Prof. Dr. Sebastian Wild

5 Integer Data Structures

- 5.1 Integer Dictionaries
- 5.2 Perfect Hashing
- 5.3 Dynamic Perfect Hashing
- 5.4 Binary Tries
- 5.5 x-Fast Tries
- 5.6 y-Fast Tries
- 5.7 Lower Bounds

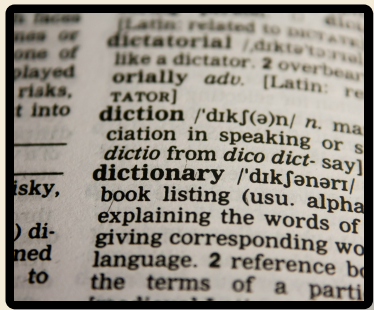
5.1 Integer Dictionaries

Recall: Symbol table ADT

Java: `java.util.Map<K,V>`, C++: `std::(unordered_)map`

Symbol table / Dictionary / Map / Associative array / key-value store:

Python dict {k:v}



- ▶ `put(k, v)` Python dict: `d[k] = v`
Put key-value pair (k, v) into table
- ▶ `get(k)` Python dict: `d[k]`
Return value associated with key k
- ▶ `delete(k)` Python dict: `del d[k]`
Remove key k (any associated value) from table
- ▶ `contains(k)` Python dict: `k in d`
Returns whether the table has a value for key k
- ▶ `isEmpty(), size(), create()`

Unlike before: Focus here on unordered case

- ▶ no rank-based queries
 - ▶ values uninteresting
- ⇒ summarize get and put as `lookup(k)`

Integer Dictionaries

For unsorted dictionaries, best solutions based on *hashing*.

- ▶ Whatever the actual objects we store, always apply a hash function first

↪ Consider here directly (dynamic) **sets of integers**

Hash Tables

Notation

- ▶ **universe** $U = [0..u)$ of allowed values
- ▶ **universe size** $u = |U|$, usually $u = 2^w$
- ▶ $S \subseteq U$ dynamic **set** of integers stored in our dictionary
- ▶ $n = |S|$ current **size** of set, $S = \{x_1, \dots, x_n\}$
- ▶ use **hash table** $T[0..m)$ with m **buckets**
- ▶ **hash function** $h : U \rightarrow [0..m)$
- ↔ element $x \in U$ stored in $T[h(x)]$
- ▶ in general: $T[h(x)]$ may be a **secondary data structure**
 - ▶ linked list ↔ *chaining hashing*
 - ▶ sequence of positions in T ↔ *open-addressing hashing*
- ▶ if $h(x) = h(y)$, we have a **collision** (between x and y)

Hashing must be randomized

Collisions are inevitable

- ▶ usually want S represented with $O(n)$ space

↪ $m = O(n)$

- ▶ typically thus $|U| = 2^w \gg m$

↪ h must have many collisions

↪ For any *fixed* hash function h , worst-case set S has many collisions

↪ cannot have meaningful worst-case time bounds (even if typical case good)

Randomized Hashing

- ▶ draw *random* $h \in \mathcal{H}$

↪ running time expected over random choice, worst-case w.r.t. S

- ▶ allowing all $\mathcal{H} = \{h \mid h : U \rightarrow [0..m)\}$ requires m^u space ⚡

↪ must choose sufficiently small \mathcal{H}

Universal Hashing

Many guarantees of varying strength studied for \mathcal{H} ... here only simplest needed

Definition 5.1 (Universal Family)

Let \mathcal{H} be a set of hash functions from U to $[m]$ and $|U| \geq m$.

Assume $h \in \mathcal{H}$ is chosen uniformly at random.

(a) Then \mathcal{H} is called a *c-universal* if

$$\forall x_1, x_2 \in U : x_1 \neq x_2 \implies \mathbb{P}[h(x_1) = h(x_2)] = \frac{c}{m}.$$

(b) \mathcal{H} is called *strongly universal* or *pairwise independent* if

$$\forall x_1, x_2 \in U, y_1, y_2 \in R : x_1 \neq x_2 \implies \mathbb{P}[h(x_1) = y_1 \wedge h(x_2) = y_2] \leq \frac{1}{m^2}. \quad \blacktriangleleft$$

Examples of universal families

Universal families

$$\blacktriangleright \mathcal{H}_1 = \{h_{ab} : a \in [1..p), b \in [0..p)\} \quad (c = 2)$$

$$\blacktriangleright \mathcal{H}_0 = \{h_{ab} : a \in [0..p), b \in [0..p)\} \quad (c = 4)$$

$$\blacktriangleright \mathcal{H}_2 = \{h_a : a \in [1..2^k), a \text{ odd}\} \quad (c = 2)$$

$$h_{ab}(x) = (a \cdot x + b \bmod p) \bmod m \quad p \text{ prime}, p \geq m$$

$$h_a(x) = (a \cdot x \bmod 2^k) \operatorname{div} 2^{k-\ell} \quad u = 2^k, m = 2^\ell$$

What does universality buy us?

Lemma 5.2 (Universal collisions)

Consider hashing n keys x_1, \dots, x_n into m buckets using a random hash function h chosen from a c -universal family of hash functions \mathcal{H} . Let C denote the number of pairwise collisions, i. e., $C = |\{(i, j) : 1 \leq i < j \leq n \wedge h(x_i) = h(x_j)\}|$.

Then $\mathbb{E}[C] < \frac{cn^2}{2m}$.

Lemma 5.3 (Universal bin size)

Consider the scenario of Lemma 5.2. Let $X_j = |T[j]|$ be the number of elements in bucket j .

(a) $\mathbb{E}[X_h(x_i)] \leq 1 + c \cdot \frac{n}{m}$

(b) Then $\mathbb{P}\left[\max X_j > \sqrt{2c} \cdot \frac{n}{\sqrt{m}}\right] \leq \frac{1}{2}$

5.2 Perfect Hashing

Perfect Hashing

A hash function $h : [u] \rightarrow [m]$ is called

- ▶ *perfect* for a set $S = \{x_1, \dots, x_n\} \subset [u]$ if $i \neq j$ implies $h(x_i) \neq h(x_j)$
- ▶ *minimal* for set $S = \{x_1, \dots, x_n\} \subset [u]$ if $m = n$

Perfect Hashing

- ▶ only possible for $n \leq m$
- ▶ stringent requirement \rightsquigarrow first focus on **static** \mathcal{X}
- ▶ further requirements
 1. Hash function must be fast to evaluate (ideally $O(1)$ time)
 2. Hash function must be small to store (ideally $O(n)$ space)
 3. should be fast to compute given \mathcal{X} (ideally $O(n)$ time)
 4. Have small m (ideally $m = \Theta(n)$)

Fredman-Komlós-Szemerédi Scheme

Theorem 5.4 (FKS Static Perfect Hashing)

Given a c -universal family of hash functions \mathcal{H} and a static set S , we can construct in $O(n)$ expected time a perfect hash function $h : S \rightarrow [m]$ with $m = O(n)$ and evaluation time $O(t)$ for t the time to evaluate functions in \mathcal{H} . ◀

Step 1: Simple, but space inefficient

Step 2: Two-tier solution

5.3 Dynamic Perfect Hashing

Dietzfelbinger et al. Scheme

Theorem 5.5 (Dynamic Perfect Hashing)

There is a randomized dictionary data structure for integers that supports $O(1)$ worst-case lookup, $O(1)$ expected amortized time insert and delete and uses $O(n)$ space for storing a set S of size n .



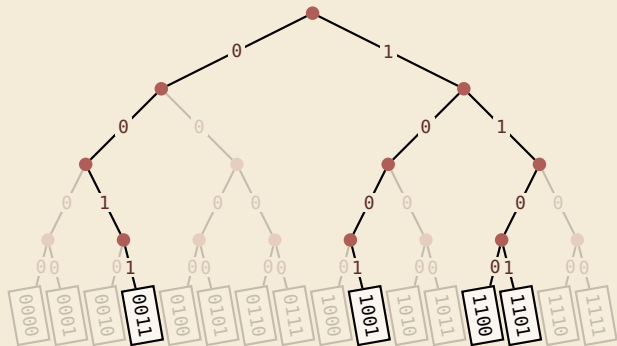
Dietzfelbinger, Karlin, Mehlhorn, Meyer auf der Heide, Rohnert, Tarjan:
Dynamic Perfect Hashing: Upper and Lower Bounds, SICOMP 1994

Dynamic Perfect Hashing – Discussion

- 👍 Strongest possible guarantee for lookup
- 👍 All operations in $O(1)$ clearly optimal (even though expected amortized)
 - ↪ (almost) perfect dictionary?
- 👎 Main downside (of hashing generally): No sorted-dictionary operations

5.4 Binary Tries

Tries



- ▶ Numbers as trie over bit sequence
- ▶ Leaves doubly linked lists in sorted order
- ▶ Each node stores *jump* pointer to extremal leaf
 - ▶ if node is left child, *jump* is leftmost leaf in subtree
 - ▶ if node is right child, *jump* is rightmost leaf

Binary Trie – Result

Theorem 5.6

A *binary trie* supports updates and predecessor / successor queries all in $O(w)$ time per operation for a set of w -bit integers, using $O(nw)$ bits of space. ◀

5.5 x-Fast Tries

Hashing!

5.6 y-Fast Tries

“Indirection” a.k.a. Don’t Change Too Fast!

5.7 Lower Bounds

The Cell-Probe Model