

Data structures II

NTIN067

Jirka Fink

<https://ktiml.mff.cuni.cz/~fink/>

Katedra teoretické informatiky a matematické logiky
Matematicko-fyzikální fakulta
Univerzita Karlova v Praze

Summer semester 2018/19
Last change March 18, 2019

License: Creative Commons BY-NC-SA 4.0

Jirka Fink Data Structures II 1

Jirka Fink: Data Structures II

General information

E-mail fink@ktiml.mff.cuni.cz

Homepage <https://ktiml.mff.cuni.cz/~fink/>

Consultations Individual schedule

Jirka Fink Data Structures II 3

Computational model Word-RAM

Description

- A word is a w -bit integer
- A memory an array of words indexed by words
- The size of memory is 2^w , so we assume that $w = \Omega(\log n)$
- Operations of words in constant time:
 - Arithmetical operations are $+$, $-$, $*$, $/$, \bmod
 - Bit-wise operations $\&$, $|$, \ll , \gg , $<<$
 - Comparisons $=$, $<$, \leq , $>$, \geq
- Other operations in constant time: (un)conditional jumps, assignments, memory accesses, etc.
- Inputs and outputs are stored in memory

Jirka Fink Data Structures II 5

Independent repeated trials

Markov inequality

If X is a non-negative random variable and $c > 1$, then $P[X < cE[X]] > \frac{c-1}{c}$.

Expected number of trial using probability

Let V be an event that occurs in a trial with probability p . The expected number of trials to first occurrence of V in a sequence of independent trials is $\frac{1}{p}$.

Expected number of trial using mean

If X is a non-negative random variable and $c > 1$. The expected number of trials to first occurrence of $X \leq cE[X]$ in a sequence of independent trials is at most $\frac{c}{c-1}$.

Example

If the expected number of collisions of a randomly chosen hashing function h is k , then the expected number of independent trials to the first occurrence of a hashing function h with at most $2k$ collisions is 2.

Jirka Fink Data Structures II 7

Content

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 Dynamization and persistence
- 5 Cache-oblivious structures
- 6 Succinct data structures
- 7 Computation in stream model
- 8 Literatura

Jirka Fink Data Structures II 2

Content

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 Dynamization and persistence
- 5 Cache-oblivious structures
- 6 Succinct data structures
- 7 Computation in stream model
- 8 Literatura

Jirka Fink Data Structures II 4

Static dictionaries

Notations

- Universe U of all elements (words)
- Store $S \subseteq U$ of size n in a data structure
- Using hashing, we store S in a table $M = [m] = \{0, \dots, m-1\}$ of size m

Goal

Create a data structure determining whether a given element of U belongs to S .

Methods

	Build	Member	
Search tree	$n \log n$	$\log n$	optimal in the comparison model
Cuckoo	n (exp.)	1	$\log n$ -independent
FKS	n (exp.)	1	2-independent
	$n \log n$	1	deterministic

Jirka Fink Data Structures II 6

Universal hashing systems

c -universal hashing system

A hashing system \mathcal{H} of functions $h: U \rightarrow M$ is c -universal for $c > 1$ if a uniformly chosen h from \mathcal{H} satisfies $P[h(x) = h(y)] \leq \frac{c}{m}$ for every $x, y \in U$ and $x \neq y$.

k -independent hashing system

A hashing system \mathcal{H} of functions $h: U \rightarrow M$ is k -independent for $k \in \mathbb{N}$ if a uniformly chosen h from \mathcal{H} satisfies $P[h(x_i) = z_i \text{ for all } i = 1, \dots, k] = O(\frac{1}{m^k})$ for all pairwise different $x_1, \dots, x_k \in U$ and all $z_1, \dots, z_k \in M$.

Example: System Multiply-mod-prime

- Let p be a prime greater than $|U|$
- $h_{a,b}(x) = (ax + b \bmod p) \bmod m$
- $\mathcal{H} = \{h_{a,b}; a, b \in [p], a \neq 0\}$
- System \mathcal{H} is 1-universal and 2-independent but it is not 3-independent

Jirka Fink Data Structures II 8

Goal

A static dictionary for n w -bit keys with constant lookup time and a space consumption of $\mathcal{O}(n)$ words can be constructed in $\mathcal{O}(n \log n)$ time on w -word-RAM. The algorithm is weakly non-uniform, i.e. requires certain precomputed constants dependent on w .

Overview

- 1 Create a function $f_1 : [2^w] \rightarrow [2^{4w}]$ which is an error-correcting code of relative minimum distance $\delta > 0$.
 - 2 Create a function $f_2 : [2^{4w}] \rightarrow [\mathcal{O}(n^k)]$ which is an injection on $f_1(S)$
 - 3 Create a function $f_3 : [\mathcal{O}(n^k)] \rightarrow [\mathcal{O}(n^2)]$ which is an injection on $f_2(f_1(S))$
 - 4 Create a function $f_4 : [\mathcal{O}(n^2)] \rightarrow [\mathcal{O}(n)]$ which is an injection on $f_3(f_2(f_1(S)))$
- $f_4 \circ f_3 \circ f_2 \circ f_1$ can be computed in constant time
 - f_2, f_3, f_4 can be found in time $\mathcal{O}(n \log n)$
 - f_1 can be precomputed in time $\mathcal{O}(w)$

- 1 The number of remaining bits is at most r .
- 2 Since $r \geq \frac{w}{2}$.

Goal

Find a function $h : U \rightarrow [2^r]$ with $U = [\mathcal{O}(n^2)]$ and $r = \max\{\frac{w}{2}, 3 + \log n\}$ s.t.

- h is perfect on $S \subseteq U$ of size n and
- h can be computed in constant time and
- space consumption $\mathcal{O}(n)$ for finding and storing h and
- h can be found in $\mathcal{O}(n \log n)$ worst case time.
 - First, expected $\mathcal{O}(n)$ time,
 - then derandomize to $\mathcal{O}(n \log n)$ worst case time.

$x \in U$ is a point $(\alpha(x), \beta(x))$ in a $(\mathcal{O}(n) \times \mathcal{O}(n))$ -table

For $x \in U$, let $\alpha(x)$ denote the first r bits of x and $\beta(x)$ denotes the remaining bits. ① Then $x \mapsto (\alpha(x), \beta(x))$ is an injection (e.i. perfect on U). ②

Static dictionaries: $[\mathcal{O}(n^2)] \rightarrow [\mathcal{O}(n)]$

Definition

For $q \geq 0$ and functions $f, g : U \rightarrow [2^r]$, the pair (f, g) is q -good if

- f has at most q collisions and
- $x \mapsto (f(x), g(x))$ is perfect on S .

The number of collisions is the number of pairs $\{x, y\} \subseteq S$ such that $f(x) = f(y)$.

Lemma

Suppose that (f, g) is q -good and $r \geq 3 + \log n$. Then, for every $v \in [2^r]$ there exists

$a_v \in [2^r]$ such that $(x \mapsto g(x) \oplus a_{f(x)}, f)$ is q' -good where $q' = \begin{cases} 0 & \text{if } q \leq n \\ n & \text{otherwise.} \end{cases}$

All values a_v can be computed in expected time $\mathcal{O}(n)$ and space $\mathcal{O}(n)$ worst case.

Application: Randomized construction of a mapping $[\mathcal{O}(n^2)] \rightarrow [\mathcal{O}(n)]$

- (α, β) is $\binom{n}{2}$ -good
- $(x \mapsto \beta(x) \oplus a_{\alpha(x)}, \alpha) =: (\alpha', \beta')$ is n -good
- $(x \mapsto \beta'(x) \oplus a_{\alpha'(x)}, \alpha') =: (\alpha'', \beta'')$ is 0-good, so α'' is perfect

Static dictionaries: $[\mathcal{O}(n^2)] \rightarrow [\mathcal{O}(n)]$

Lemma

Suppose that (f, g) is q -good and $r \geq 3 + \log n$. Then, for every $v \in [2^r]$ there exists

$a_v \in [2^r]$ such that $(x \mapsto g(x) \oplus a_{f(x)}, f)$ is q' -good where $q' = \begin{cases} 0 & \text{if } q \leq n \\ n & \text{otherwise.} \end{cases}$

All values a_v can be computed in expected time $\mathcal{O}(n)$ and space $\mathcal{O}(n)$ worst case.

Proof ($q' \leq n$)

- 1 Let $h(x) = g(x) \oplus a_{f(x)}$
- 2 If $x, y \in S$ and $x \neq y$ and $f(x) = f(y)$, then $g(x) \neq g(y)$ and $h(x) \neq h(y)$ ①
- 3 If $f(x) \neq f(y)$, then $P[h(x) = h(y)] = \frac{1}{2^r}$ where $a_v \sim U[2^r]$ independently for all $v \in [2^r]$ ②
- 4 $E[|\{x, y\} \subseteq S; h(x) = h(y)\}|] \leq \binom{n}{2} / 2^r < \frac{n}{16}$ ③
- 5 The expected number of trials to generate h with at most n collisions is $\mathcal{O}(1)$.

Static dictionaries: $[\mathcal{O}(n^2)] \rightarrow [\mathcal{O}(n)]$

Lemma

Suppose that (f, g) is q -good and $r \geq 3 + \log n$. Then, for every $v \in [2^r]$ there exists

$a_v \in [2^r]$ such that $(x \mapsto g(x) \oplus a_{f(x)}, f)$ is q' -good where $q' = \begin{cases} 0 & \text{if } q \leq n \\ n & \text{otherwise.} \end{cases}$

All values a_v can be computed in expected time $\mathcal{O}(n)$ and space $\mathcal{O}(n)$ worst case.

Proof ($q \leq n$ implies $q' = 0$)

- 1 Let $S_v = \{x \in S; f(x) = v\}$
- 2 Order S_v by non-increasing size, i.e. $|S_{v_1}| \geq |S_{v_2}| \geq \dots \geq |S_{v_{2^r}}|$
- 3 For $j = 1, \dots, 2^r$ we find a_{v_j} such that h is perfect ①
- 4 For $a_{v_j} \sim U[2^r]$ it holds $E[|\{(x, y) \in S_{v_j} \times S_{v_j}; h(x) = h(y)\}|] \leq |S_{v_j}| |S_{v_j}| P[h(x) = h(y)]$ ②
 - $\leq \sum_{i=1}^{j-1} |S_{v_j}| |S_{v_i}| / 2^r$ ③
 - $\leq \sum_{i=1}^{j-1} |S_{v_i}^2| / 2^r$ ④
 - $\leq \sum_{i=1}^{j-1} \binom{|S_{v_i}|}{2} / 2^{r-2}$ ⑤
 - $\leq q / 2^{r-2} \leq \frac{1}{2}$
- 5 The expected number of trials to generate a_{v_j} such that h has no collision is $\mathcal{O}(1)$. ⑥ ⑦

Definition

For $q \geq 0$ and functions $f, g : U \rightarrow [2^r]$, the pair (f, g) is q -good if

- f has at most q collisions and
- $x \mapsto (f(x), g(x))$ is perfect on S .

The number of collisions is the number of pairs $\{x, y\} \subseteq S$ such that $f(x) = f(y)$.

Lemma

Suppose that (f, g) is q -good and $r \geq 3 + \log n$. Then, for every $v \in [2^r]$ there exists

$a_v \in [2^r]$ such that $(x \mapsto g(x) \oplus a_{f(x)}, f)$ is q' -good where $q' = \begin{cases} 0 & \text{if } q \leq n \\ n & \text{otherwise.} \end{cases}$

All values a_v can be computed in expected time $\mathcal{O}(n)$ and space $\mathcal{O}(n)$ worst case.

Application: Randomized construction of a mapping $[\mathcal{O}(n^2)] \rightarrow [\mathcal{O}(n)]$

- (α, β) is $\binom{n}{2}$ -good
- $(x \mapsto \beta(x) \oplus a_{\alpha(x)}, \alpha) =: (\alpha', \beta')$ is n -good
- $(x \mapsto \beta'(x) \oplus a_{\alpha'(x)}, \alpha') =: (\alpha'', \beta'')$ is 0-good, so α'' is perfect

- 1 Since $x \mapsto (f(x), g(x))$ is perfect on S , $g(x) \neq g(y)$. From $g(x) \oplus a_{f(x)} = g(y) \oplus a_{f(y)} \neq g(y) \oplus a_{f(y)}$ it follows that $h(x) \neq h(y)$.
- 2 For every $v \in [2^r]$ we randomly and independently choose a_v from the uniform distribution on $[2^r]$. Then,

$$\begin{aligned} h(x) &= h(y) \\ g(x) \oplus a_{f(x)} &= g(y) \oplus a_{f(y)} \\ a_{f(x)} &= g(x) \oplus g(y) \oplus a_{f(y)} \end{aligned}$$

Since $([2^r], \oplus)$ is an Abelian group, $b \mapsto b \oplus c$ is a bijection on $[2^r]$ for every $c \in [2^r]$ and so $a_{f(y)} \sim U[2^r]$, it follows that $g(x) \oplus g(y) \oplus a_{f(y)} \sim U[2^r]$. Since $a_{f(x)}$ and $a_{f(y)}$ are independent, also $a_{f(x)}$ and $g(x) \oplus g(y) \oplus a_{f(y)}$ are independent. Hence, $P[h(x) = h(y)] = \frac{1}{2^r}$.

- 3 Use the linearity of expectation and substitute r .

- 1 Note that we must find h without collisions. To be precise, we iteratively find a_{v_j} for j from 1 to 2^r such that it holds $h(x) \neq h(y)$ for every $x \in S_{v_j}$ and $y \in S_{<j}$ where $S_{<j} = \bigcup_{i=1}^{j-1} S_{v_i}$.

- 2 Linearity of expectation

- 3 Definition of $S_{<j}$

- 4 $|S_{v_j}| \geq |S_{v_j}|$

- 5 From this point, assume that $|S_{v_j}| \geq 2$.

- 6 In order to verify that h has no collision, we use a counter $m_v = |\{y \in S_{<j}; h(y) = v\}|$. For every j we can count the collisions and update m_v in time $\mathcal{O}(|S_j|)$. The expected time to find all $a_{v_j} = \sum_j \mathcal{O}(|S_j|) = \mathcal{O}(n)$.

- 7 For $S_{v_j} = \{x\}$ we can find v with $m_v = 0$ and set $a_{v_j} = v \oplus g(x)$.

Static dictionaries: $[O(n^2)] \rightarrow [O(n)]$

Derandomization

Let $L_k(a) = |\{(x, y) \in S_{y_j} \times S_{<j}; (g(x) \oplus a)_{[k]} = (h(y))_{[k]}\}|$ \odot
 $(a)_i$ denotes the i -th bit of a and $(a)_M$ denotes the vector of all bits $(a)_i$ for $i \in M \subseteq [r]$ \odot

```

1 for  $j \leftarrow 1$  to  $2^r$  do
2    $a_{y_j} \leftarrow 0$ 
3   for  $k \leftarrow 0$  to  $r - 1$  do
4     if  $L_{k+1}(a_{y_j}) > L_{k+1}(a_{y_j} + 2^k)$  then
5        $a_{y_j} \leftarrow a_{y_j} + 2^k$ 

```

Proof (goodness of (h, f))

- $L_k(a) + L_k(a \oplus 2^k) = L_{k-1}(a)$ for every $a \in [2^k]$ and $k \in [r]$
- $L_k(a_{y_j}) \leq \frac{L_{k-1}(a_{y_j})}{2} \leq \frac{L_0(a_{y_j})}{2^k} = \frac{|S_{y_j}| |S_{<j}|}{2^k}$
- The total number of collision is at most $\sum_j L_r(a_{y_j}) \leq \sum_j 2^{-r} |S_{y_j}| |S_{<j}| \leq \sum_{i < j} 2^{-r} |S_{y_j}| |S_{y_i}| \leq 2^{-r-1} (\sum_i |S_{y_i}|)^2 < \frac{n}{16}$
- If $q \leq n$, then the number of collision with S_{y_j} is $L_r(a_{y_j}) \leq \sum_{i < j} 2^{-r} |S_{y_j}| |S_{y_i}| \leq \sum_{i < j} 2^{2-r} \binom{S_{y_j}}{2} \leq 2^{2-r} q \leq \frac{1}{2}$

Jirka Fink

Data Structures II

14

Static dictionaries: $[O(n^2)] \rightarrow [O(n)]$

Derandomization

Let $L_k(a) = |\{(x, y) \in S_{y_j} \times S_{<j}; (g(x) \oplus a)_{[k]} = (h(y))_{[k]}\}|$
 $(a)_i$ denotes the i -th bit of a and $(a)_M$ denotes the vector of all bits $(a)_i$ for $i \in M \subseteq [r]$

```

1 for  $j \leftarrow 1$  to  $2^r$  do
2    $a_{y_j} \leftarrow 0$ 
3   for  $k \leftarrow 0$  to  $r - 1$  do
4     if  $L_{k+1}(a_{y_j}) > L_{k+1}(a_{y_j} + 2^k)$  then
5        $a_{y_j} \leftarrow a_{y_j} + 2^k$ 

```

Proof (Complexity)

- In order to compute $L_k(a)$, we build a binary tree (trie)
- Every vertex $a \in [2^k]$ of the k -th level has a counter $c_k(a) = |\{y \in S_{<j}; (h(y))_{[k]} = (a)_{[k]}\}|$
- $L_k(a) = \sum_{x \in S_{y_j}} c_k(g(x) \oplus a)$ can be computed in $O(|S_{y_j}|)$ time
- After the j -th step, counters can be updated in $O(|S_{y_j}|r)$ time
- Total time is $\sum_j |S_{y_j}|r = O(n \log n)$

Jirka Fink

Data Structures II

15

Static dictionaries: Error-correcting code

Definition

- The Hamming distance between $x \in [2^w]$ and $y \in [2^w]$ is the number of bits in which x and y differ.
- $\psi: [2^w] \rightarrow [2^{4w}]$ is an error correcting code of relative minimum distance $\delta > 0$ if the Hamming distance between $\psi(x)$ and $\psi(y)$ is at least $4w\delta$ for every distinct $x, y \in [2^w]$.

Lemma

Let \mathcal{H} be a 2-universal hashing system of function $[2^w] \rightarrow [2^{4w}]$. For every δ with $1/4w < \delta \leq 1/2$, the probability that $h \sim \mathcal{H}$ is an error correcting code of relative minimum distance $\delta > 0$ is at least $1 - ((\frac{e}{\delta})^{4\delta}/4)^w$. \odot

Proof

- For $x \in [2^{4w}]$ the number of y within Hamming distance k is at most $(\frac{4ew}{k})^k$. \odot
- For $x \neq y$, $P(\text{Hamming distance between } x \text{ and } y \leq k) \leq 2^{1-4w} (\frac{4ew}{k})^k$
- The probability that this happens for any of the $\binom{2^w}{2} < 2^{2w-1}$ pairs is at most $((\frac{e}{\delta})^{4\delta}/4)^w$ \odot

Jirka Fink

Data Structures II

17

Static dictionaries: $[2^w] \rightarrow [O(n^k)]$

Lemma

Let $\psi: [2^w] \rightarrow [2^{4w}]$ be an error correcting code of relative minimum distance $\delta > 0$ and $S \subseteq U = [2^w]$ of size n . There exists a set $D \subseteq [4w]$ with $|D| \leq 2 \log n / \log \frac{1}{1-\delta}$ such that for every pair x, y of distinct elements of S it holds $(\psi(x))_D \neq (\psi(y))_D$.

Proof

- For $D \subseteq [4w]$ and $v \in [2^{|D|}]$ let $C(D, v) = \{x \in S; (\psi(x))_D = v\}$ \odot
- The set of colliding pairs of D is $B(D) = \bigcup_{v \in [2^{|D|}]} \binom{C(D, v)}{2}$
- We construct $D_0 \subseteq D_1 \subseteq \dots \subseteq D_k$ such that $|D_i| = i$ and $|B(D_i)| < (1 - \delta)^i n^2 / 2$ \odot
- Let $l(d) = \{x, y \in B(D_i); (\psi(x))_d = (\psi(y))_d\}$ be the colliding pairs indistinguishable by $d \in [4w] \setminus D_i$
- Let $l = \sum_{d \in [4w] \setminus D_i} |l(d)|$
- Every pair $\{x, y\} \in B(D_i)$ contributes to l by at most $4w - i - 4w\delta < 4w(1 - \delta)$, so $l \leq 4w(1 - \delta) |B(D_i)|$
- By averaging, there exists $d \in [4w] \setminus D_i$ such that $|l(d)| \leq \frac{l}{4w-i} \leq (1 - \delta) |B(D_i)|$ \odot
- Let $D_{i+1} = D_i \cup \{d\}$. Hence, $|B(D_{i+1})| = |l(d)| \leq (1 - \delta) |B(D_i)|$
- By setting $k = \left\lceil 2 \log n / \log \frac{1}{1-\delta} \right\rceil$ we obtain $|B(D_k)| < 1$.

Jirka Fink

Data Structures II

18

- Where $k \in [r]$ and $a \in [2^k]$

- Our goal is to iteratively and deterministically compute a_{y_j} for j from 1 to 2^r . The value of a_{y_k} is computed by bits from the least significant to the most significant bit. $L_k(a)$ determines the number of collision between S_{y_j} and $S_{<j}$ if we consider only last k bits.

Jirka Fink

Data Structures II

14

Static dictionaries: $[O(n^k)] \rightarrow [O(n^2)]$

Approach

- Every $x \in U = [O(n^k)]$ can be regarded as constant-length string over an alphabet of size n
- Build n -way branching compressed trie of string S
- The number of leaves is $|S| = n$, so the total number of vertices is at most $2n$
- Build static $[O(n^2)] \rightarrow [O(n)]$ dictionary for pairs (vertex of the trie, letter) which returns a child of the vertex
- One polynomial-size-universe lookup is evaluated using a constant number of quadratic-size-universe lookups
- Space complexity is $O(n)$ and this dictionary is constructed in $O(n \log n)$ time

Jirka Fink

Data Structures II

16

- For $\delta < \frac{1}{4w}$ it holds that $4w\delta < 1$ and the identity is an error correcting code of relative minimum distance δ .
- The number of $y \in [2^{4w}]$ within Hamming distance $k \geq 1$ from a fixed $x \in [2^{4w}]$ is $\sum_{i=0}^k \binom{4w}{i} \leq (\frac{4w}{k})^k \sum_{i=0}^k \binom{4w}{i} (\frac{k}{4w})^i \leq (\frac{4w}{k})^k (1 + \frac{k}{4w})^{4w} \leq (\frac{4w}{k})^k e^k \leq (\frac{4ew}{k})^k$ using the binomial theorem.
- By setting $k = \lceil 4w\delta \rceil$ we obtain $2^{2w-1} 2^{1-4w} (\frac{4ew}{k})^k \leq 2^{-2w} (\frac{4ew}{4w\delta})^{4w\delta} = (2^{-2} (\frac{e}{\delta})^{4\delta})^w$

Jirka Fink

Data Structures II

17

- Note that for every $D \subseteq [4w]$ the set S is split into $2^{|D|}$ disjoint clusters $C(D, v)$ for $v \in [2^{|D|}]$.
- For $i = 0$ it holds that $D_0 = \emptyset$ and $B(D_0) = \binom{n}{2} < \frac{n^2}{2}$.
- A bit $d \in [4w] \setminus D_i$ with $|l(d)| \leq (1 - \delta) |B(D_i)|$ can be found in $O(wn)$ time as follows. We list all clusters $C(D_i, v)$ of size at least two. Every cluster has a list of all elements. So, $l(d)$ for one $d \in [4w] \setminus D_i$ can be determined in $O(n)$ time and we can process all d in $O(wn)$ time. Then, all lists can be updated in $O(n)$ time. Using word-level parallelism, the time complexity can be improved to $O(n)$.

Jirka Fink

Data Structures II

18

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 Dynamization and persistence
- 5 Cache-oblivious structures
- 6 Succinct data structures
- 7 Computation in stream model
- 8 Literatura

- 1 This theorem is similar to the linear speed theorem for Turing machines.
- 2 The assumption $S(n) = o(2^w)$ is needed to ensure a sufficient amount of memory for linear data decompression.
- 3 Bit-wise operations and comparisons are trivial. In order to simulate arithmetical operations, every $\mathcal{O}(w)$ bits word is split into $\frac{w}{2}$ bits words to handle integer overflows and carry.

- 1 Reading initialized memory must be permitted but an arbitrary value can be read.

van Emde Boas tree [3]

Notation

For a t -bit integer $x \in U$ let

- $r(x)$ be the left (upper) $\lceil t/2 \rceil$ bits of x and
- $l(x)$ be the left (lower) $\lceil t/2 \rceil$ bits of x .

For a set of elements S let

- $S_i = \{r(x) \in S \mid \min S; l(x) = i\}$ for $i \in l(U)$
- $P = \{S_i; S_i \neq \emptyset\}$

A node of van Emde Boas tree contains

- Pointers to the minimum and maximum of S
- An array of vEB trees for all sets S_i
- A primary vEB tree of elements P for finding non-empty buckets

Space complexity

- Terrible: $\mathcal{O}(U)$, but we will improve it later.
- Requires initialized memory

Theorem (Linear compression and speed up)

An algorithm running on a $\mathcal{O}(w)$ -word-RAM with time complexity $T(n)$ and space complexity $S(n)$ can be simulated on w -word-RAM with time complexity $\mathcal{O}(T(n))$ and space complexity $\mathcal{O}(S(n))$ assuming $S(n) = o(2^w)$. ① ②

Proof

- Every word of $\mathcal{O}(w)$ -word-RAM is replaced by $\mathcal{O}(1)$ words on w -word-RAM, so the space complexity is $\mathcal{O}(T(n))$.
- Every instruction on $\mathcal{O}(w)$ -word-RAM is replaced by $\mathcal{O}(1)$ instruction on w -word-RAM. ③

RAM with uninitialized memory

Theorem

Every program with time complexity $T(n)$ and space complexity $S(n)$ which assumes initialized memory by zeros can be simulated by a program with time complexity $\mathcal{O}(T(n))$ and space complexity $S(n)$ which do not assume initialized memory. ①

Proof

Simulation uses the following variables:

- M: An array of the original memory
- N: The number of already initialized cells by the original program
- I: An array of indices of all initialized cells
- X: A reverse table of I, i.e. if j is initialized, then $I[X[j]] = j$

A read cell j is initialized if and only if $X[j] < N$ and $I[X[j]] = j$.

Integer data structures

Goal

Create a dynamic data structure storing a set S of intergers from a universe U with is can find the predecessor and the successor of a element $x \in U$.

- Predecessor: $\text{PRED}(x \in U) = \max \{y \in S; y < x\}$
- Successor $\text{SUCC}(x \in U) = \min \{y \in S; y > x\}$

An invalid element is returned if no such element exists.

Our knowledge

$\mathcal{O}(\log n)$ query and update using search trees

van Emde Boas tree: Operation FIND(x)

Algorithm

```

1 Procedure FIND ( $X, S$ )
2   if  $S = \emptyset$  then
3     return None
4   if  $x = \min(S)$  then
5     return  $\min(S)$ 
6   return FIND ( $r(x), S_{l(x)}$ )

```

Analysis

- The height of the tree is $\mathcal{O}(\log \log U)$
- Time complexity is $\mathcal{O}(\log \log U)$

Algorithm

```

1 Procedure Succ (x, S)
2   if S =  $\emptyset$  or x  $\geq$  max(S) then
3     return None
4   if x < min(S) then
5     return min(S)
6   i  $\leftarrow$  l(x)
7   if Si  $\neq \emptyset$  and x < max(Si) then
8     return Succ (r(x), Si)
9   i  $\leftarrow$  Succ (i, P)
10  return min(Si)

```

Analysis

- The recursion is called at most once for every level of the tree
- Time complexity is $\mathcal{O}(\log \log U)$

Jirka Fink

Data Structures II

25

van Emde Boas tree: Operation DELETE(x)

Algorithm ①

```

1 Procedure Delete (x, S)
2   if S =  $\emptyset$  or x = min(S) = max(S) then
3     S  $\leftarrow \emptyset$  and return
4   if x = min(S) then
5     min(S)  $\leftarrow$  DeleteMin (S) and return
6   i  $\leftarrow$  l(x)
7   if x = min(Si) = max(Si) then
8     Delete (i, P)
9   Delete (r(x), Si)
10  max(S)  $\leftarrow$  max(Smax(P))
11 Procedure DeleteMin (S) ②
12   i  $\leftarrow$  min(P)
13   if min(Si) < max(Si) then
14     return DeleteMin (Si)
15   Si  $\leftarrow \emptyset$ 
16   i  $\leftarrow$  DeleteMin (P)
17   return min(Si)

```

Jirka Fink

Data Structures II

27

x-fast trie (Willard [4])

First step

- A set *S* is represented using an array *A* such that *x* is stored at *A*[*x*]
- Non-empty positions of *A* are interconnected using a linked list
- Build a binary tree *T* whose leaves are cells of *A*
- Every node *v* of *T* stores the minimum and maximum of the subtree of *v*

Operation Succ(x)

- If *A*[*x*] is non-empty, then return *A*[*x*] \rightarrow *next*
- Find the largest empty subtree *v* containing *x* and its parent *p*
- If *v* is a left child of *p*, then return *p* \rightarrow *min*
- Return *p* \rightarrow *max* \rightarrow *next*

Complexity

- Space: $\mathcal{O}(U)$
- Succ: $\mathcal{O}(\log U)$
- INSERT, DELETE: $\mathcal{O}(\log U)$

Jirka Fink

Data Structures II

28

y-fast trie (Willard [4])

Structure

- Split the ordered set *S* into $\Theta(n/B)$ blocks *S_i* of size $\Theta(B)$ where $B = \log U$
- For every block *S_i* choose a representative *r_i* and let *R* be the set of all representatives
- Create x-fast trie for representatives *R*, called a primary tree
- For every block *S_i*, create a balanced search tree *T_i*, called block trees

Space complexity

- Every elements of *S* is stored in one node of one block tree
- Space complexity of all block trees is $\mathcal{O}(n)$
- Space complexity of the primary tree is $\mathcal{O}(|R|w) = \mathcal{O}(n)$ since $|R| = \Theta(n/\log U)$ and $|w| = \Theta(\log U)$

Jirka Fink

Data Structures II

30

Algorithm

```

1 Procedure Insert (x, S)
2   if S =  $\emptyset$  then
3     min(S)  $\leftarrow$  max(S)  $\leftarrow$  x and return
4   if x < min(S) then
5     swap(x, min(S))
6   if x > max(S) then
7     max(S)  $\leftarrow$  x
8   i  $\leftarrow$  l(x)
9   if Si =  $\emptyset$  then
10    INSERT (i, P)
11  INSERT (r(x), Si)

```

Analysis

- If *S_i* = \emptyset , then *INSERT* (*r*(*x*), *S_i*) only sets *min*(*S_i*) and *max*(*S_i*)
- A non-trivial recursion is called at most once for every level of the tree
- Time complexity is $\mathcal{O}(\log \log U)$

Jirka Fink

Data Structures II

26

- Time complexity is also $\mathcal{O}(\log \log U)$
- Function *DELETETMIN* also returns the new minimal value.

x-fast trie: Improvements

Second step: Speed up

- The largest empty subtree *v* containing *x* can be found using binary search
- Time complexity of Succis $\mathcal{O}(\log \log U)$

Third step: Space reduction

- The number of nodes of *T* with non-empty subtree is at most *nw*
- Store them in a hash table (e.g. Cuckoo)

x-fast trie

- Space complexity $\mathcal{O}(nw)$
- Succ: $\mathcal{O}(\log \log U)$ worst case
- INSERT, DELETE: $\mathcal{O}(\log U)$ expected amortized

Jirka Fink

Data Structures II

29

y-fast trie: Operations FIND and Succ

Operation FIND(x)

- Find the predecessor *r_i* and the successor *r_j* representatives of *x* in the primary tree
- Element *x* lives in *T_i* or *T_j*, so search both block trees ①

Time complexity

- Representative are found in time $\mathcal{O}(\log \log U)$
- Searching block trees takes $\mathcal{O}(\log |B|) = \mathcal{O}(\log \log U)$
- Time complexity of operation FINDis $\mathcal{O}(\log \log U)$

Operation Succ(x)

- If *x* $\notin S$, then FIND(*x*) finds both the predecessor and the successor of *x*
- If *x* $\in S$, we use a linked list to find the successor
- Time complexity of operation Succis $\mathcal{O}(\log \log U)$

Jirka Fink

Data Structures II

31

- If x is a representative, we search only one tree.

- Remind that x-fast trie uses Cuckoo hashing

RAM as a vector computer

Replicate(a)

- Create n copies of a
- $\text{Replicate}(a) = a \cdot (0^b 1)^n$

Sum(x)

- Return the sum $\sum_{i=0}^{n-1} x_i$
- $\text{Sum}(x) = x \bmod 2^{b+1} - 1$
- $x = \sum_{i=0}^{n-1} x_i 2^{bi+1} \equiv \sum_{i=0}^{n-1} x_i \bmod 2^{b+1} - 1$

Cmp(x, y)

- return z such that $z_i = 1$ if $x_i \geq y_i$ and 0 otherwise
- in the difference $(1x_i) - (0y_i)$, the first bit is 1 if and only if $x_i \geq y_i$
- $\text{Cmp}(x, y) = (((x|(10^b)^n) - y) \gg b) \& (0^b 1)^n$

Min(x, y)

- $m = \text{Cmp}(x, y) \cdot 1^b$
- m filters positions of x where $x_i \geq y_i$
- $\text{Min}(x, y) = (x \& m) | (y \& \bar{m})$

RAM as a vector computer

Unpack(a)

- Create a vector x such that x_i equals to the i -th bit a_i of a
- $y = \text{Replicate}(a) \& (2^{b-1}, \dots, 1)$, i.e. $y_i = 2^i a_i$
- $\text{Unpack}(a) = \text{Cmp}(y, (0^b 1)^n)$
- $\text{Unpack}_\pi(a)$ according to a permutation π

Pack(a)

- Consider that b is decreased by one and apply Sum
- Since the vector is not properly formatted, modulo cannot be used
- $\text{Pack}(a) = ((a \cdot (0^{b-1} 1)^n) \gg bn) \& 1^b$

y-fast trie: Operations INSERT and DELETE

Without balancing sizes of blocks

- Find the proper block tree: $\mathcal{O}(\log \log U)$ time
- Update the block tree: $\mathcal{O}(\log B) = \mathcal{O}(\log \log U)$ time
- However, the size of blocks must be $\mathcal{O}(\log B)$ and number of blocks must be $\mathcal{O}(n/\log B)$

Balancing sizes of blocks

- If a block is too large, split it
- If a block is too small, merge it with a sibling or move $\Omega(B)$ elements from a sibling
- Updates blocks trees takes $\mathcal{O}(B)$ time
- Updates the primary tree takes $\mathcal{O}(\log U)$ time
- Homework: Create rules which ensures that balancing occur once every $\Omega(B)$ updates
- Amortized cost of balancing is $\mathcal{O}(1)$
- Time complexity is $\mathcal{O}(\log \log U)$ expected amortized ①

RAM as a vector computer

Representation of a vector

- Consider b -bits integers x_0, \dots, x_{n-1}
- Every integers is prepended by 0 and concatenated, i.e. a vector is represented as $\sum_{i=0}^n x_i 2^{i(b+1)}$
- Using linear compression theorem, it suffices to assume that $nb = \mathcal{O}(w)$

Operation Read(x, i)

- Returns the i -th integers from a vector x
- $\text{Read}(x, i) = x \gg 2^{i(b+1)} \& 1^b$
- where 1^b is a binary number with b ones, i.e. $2^{b+1} - 1$

Operation Write(x, i, a)

- Write a into the i -th position of x
- $\text{Write}(x, i, a) = (x \& M_i) | (a \ll 2^{i(b+1)})$
- where $M_i = \uparrow^{(b+1)(n-i-1)} 0^{b+1} \uparrow^{(b+1)i}$ is a mask cleaning the i -th position

RAM as a vector computer

Rank(x, a)

- The number of elements of x smaller than a
- $\text{Rank}(x, a) = \text{Sum}(\text{Cmp}(x, \text{Replicate}(a)))$

Insert(x, a)

- Insert a into a sorted vector x
- $i = \text{Rank}(x, a)$ is the position for a in x
- Filter and move x_i, \dots, x_{n-1} to left

RAM as a vector computer assuming $w \geq \Omega(b^2)$

Weight(a)

- The Hamming weight of a
- $\text{Weight}(a) = \text{Sum}(\text{Unpack}(a))$

Permute $_\pi$ (a)

- Permute bits in a according to a permutation π
- $\text{Permute}_\pi(a) = \text{Pack}(\text{Unpack}_\pi(a))$

LSB(a)

- The least significant bit in a
- $\text{LSB}(a) = \text{Weight}(a \oplus (a - 1)) - 1$

MSB(a)

- $\text{MSB}(a) = b - 1 - \text{LSB}(\text{Permute}_{\text{revers}}(a))$

MSB(a) on $\mathcal{O}(w)$ bits

- $b = \lfloor \sqrt{w} \rfloor$
- Split a into l block of b bits and padding
- $x = (a \& (01^b)^l) | ((a \& (10^b)^l) \gg b)$
- Now, padding is zero and $x_i = 0$ if and only if i -th block is zero
- $y = \text{Cmp}(x, (0^b 1)^n)$
- $y_i = 1$ if and only if i -th block is non-zero
- $j = \text{MSB}(\text{Pack}(y))$
- j is the first non-empty block
- $z = (a \gg (b+1)j) \& 1^{b+1}$
- z is the content of the j -th block
- $\text{MSB}(a) = (b+1)j + \text{MSB}(z)$

Fusion trees

Fusion node

- Consider a trie on keys x_0, \dots, x_{k-1}
- The number of leaves is k , so the number of splitting nodes is $k-1$ ①
- Hence, at most $k-1$ levels has a splitting node ②
- Let $s(x)$ select all splitting bits from the binary code of x ③
- Observe that $s(x_0) < s(x_1) < \dots < s(x_{k-1})$
- $S = (s(x_0), \dots, s(x_{k-1}))$
- The number of bits of S is at most $k^2 = \mathcal{O}(w^{2/5})$ so S can be stored in a single w -word
- $\text{Rank}(x)$ can be computed using $\text{Rank}(s(x), S)$ in $\mathcal{O}(1)$
- However, we need to compute $s(x)$ in $\mathcal{O}(1)$
- We find $a(x)$ of length $\mathcal{O}(k^4)$ containing $s(x)$ with extra zeros on the same position of all x

Fusion trees

Lemma

For every $b_0 < \dots < b_{k-1}$ there exist $m_0 < \dots < m_{k-1}$ such that

- $b_i + m_j$ are pair-wise different for all i and j
- $b_0 + m_0 < \dots < b_{k-1} + m_{k-1}$
- $(b_{k-1} + m_{k-1}) - (b_0 + m_0) = \mathcal{O}(k^4)$

Splitting bits with extra zeros

- $y = x \& \sum_i 2^{b_i}$ filters all splitting bits
- $z = y \cdot \sum_j 2^{m_j} = \sum_i \sum_j x_{b_i} 2^{b_i+m_j}$
- There is no carry since $b_i + m_j$ are pair-wise different
- $a(x) = z \& (\sum_i 2^{b_i+m_i}) \gg (b_0 + m_0)$

Path queries in trees

Goals

- Store a tree with weights on vertices
- Query: Find the minimum on the path between given two vertices
- Local update: Change the weight on a given vertex
- Path update: Increase weights on all vertices on the path between given two vertices
- Cut: Split the tree by removing a given two vertices
- Link: Join two trees by adding an edge between two given vertices
- Evert: Change the root of a tree and the orientation of all edges on the path between old and new vertices

Fusion node

- Let $k = \mathcal{O}(w^{1/5})$
- A fusion node stores keys x_0, \dots, x_{k-1}
- A fusion node can find Rank, SUCC and PRED in $\mathcal{O}(1)$
- Construction time is Poly(k)

Fusion tree

- B -tree for $B = k$
- Nodes are Fusion nodes
- The depth of the tree is $\mathcal{O}(\log_B(n)) = \mathcal{O}\left(\frac{\log n}{\log w}\right)$
- Time complexity of Rank, SUCC and PRED is $\mathcal{O}\left(\frac{\log n}{\log w}\right)$

- 1 Splitting node is a node of degree 2.
- 2 We call them splitting levels.
- 3 Note that every level of a trie corresponds to one bit in the binary code. Bits corresponding to splitting levels are called splitting bits. Since only splitting bits are significant, $s(x)$ select these splitting bits from x .

Content

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 Dynamization and persistence
- 5 Cache-oblivious structures
- 6 Succinct data structures
- 7 Computation in stream model
- 8 Literatura

Queries in sequences

Static representation

- Store a sequence of weights w_1, \dots, w_n in an array
- Build a balanced binary search tree with leaves corresponding to the weights
- Every vertex stores the minimum of all weights in its subtree

Query: Find the minimum in a given subsequence w_i, \dots, w_j

- Process $\mathcal{O}(\log n)$ nodes on paths from w_i and w_j to the root
- Time complexity is $\mathcal{O}(\log n)$

Local update of a given weight w_i

- Update the minimum in all vertices on the path from w_i to the root
- Time complexity is $\mathcal{O}(\log n)$

Increase all weights in a subsequence w_i, \dots, w_j

- Lazy evaluation: Vertex can store an increment of all weight in its subtree
- Invariant: The actual weight w_k is the sum of the value stored in k -th leaf and all increments on the path to the root
- Increment propagation: The increment stored in a vertex u can be added to increments in children and cleared in u
- All operations propagate increment on accessed vertices
- Subsequence increase: Set the increase in nodes on paths from w_i and w_j to the root
- Time complexity is $\mathcal{O}(\log n)$

Heavy-light decomposition of rooted trees

Representation

Every vertex v stores:

- Its parent u and the list of its children
- Heavy/light flag of the edge uv
- Minimum of weights in the subtree
- First vertex of the heavy path containing v
- Order of v on the heavy path containing v
- Every heavy path has a DS for sequences
- In static version: Minimum of the prefix of heavy path to v

Link-Cut trees (Sleator, Tarjan)

Definition

- Trees are rooted
- Every edges is thick or thin
- Every vertex u has at most one child v such that uv is thick
- Thick edges form a set paths
- Every vertex lies on a thick path since trivial paths are allowed

Operations

- Expose(v): Makes the path from v to the root thick and all children of v thin
- Find the parent of v and the root of v
- Cut: Split the tree by removing a given two vertices
- Link: Join two trees by adding an edge between two given vertices
- Evert: Change the root of a tree and the orientation of all edges on the path between old and new vertices
- Find the weight of a path
- Find the minimum of all weight on a path
- Change the weight of vertex
- Increase all weights of vertices on a path

Content

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 Dynamization and persistence
- 5 **Cache-oblivious structures**
- 6 Succinct data structures
- 7 Computation in stream model
- 8 Literatura

Heavy and light edges

- Let $s(v)$ be the number of vertices in the subtree of v (including v)
- Let u be a vertex and v its child
- The edge uv is called *heavy* if $s(u) \geq s(v)/2$; and *light* otherwise

Observations

- Every vertex u has at most one child v such that uv is heavy
- Heavy edges form a set paths
- Every vertex lies on a heavy path since trivial paths are allowed
- For every vertex u the path from the root to u contains at most $\log n$ light edges and most $\log n$ prefixes of heavy paths

Heavy-light decomposition of rooted trees

Lowest Common Ancestor (LCA)

For given two vertices u and v find the last common vertex on path from the root to u and v

- On the path from u to the root, mark all heavy paths
- Find the first marked heavy path p from v to the root
- Using the order of vertices on p , determine the LCA
- Time complexity is $\mathcal{O}(\log n)$

Path query

Find the minimal weight on the path between given two vertices u and v

- Find $z = \text{LCA}(u, v)$
- The path $u \dots z \dots v$ contains:
 - $\mathcal{O}(\log n)$ light edges
 - $\mathcal{O}(\log n)$ prefixes of heavy paths
 - At most one subpath of a heavy path
- Time complexity
 - Static version: $\mathcal{O}(\log n)$
 - Dynamic version: $\mathcal{O}(\log^2 n)$

Content

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 **Dynamization and persistence**
- 5 Cache-oblivious structures
- 6 Succinct data structures
- 7 Computation in stream model
- 8 Literatura

Content

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 Dynamization and persistence
- 5 Cache-oblivious structures
- 6 **Succinct data structures**
- 7 Computation in stream model
- 8 Literatura

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 Dynamization and persistence
- 5 Cache-oblivious structures
- 6 Succinct data structures
- 7 Computation in stream model
- 8 Literatura

[1] Michael L Fredman and Dan E Willard.
Surpassing the information theoretic bound with fusion trees.
Journal of computer and system sciences, 47(3):424–436, 1993.

[2] Torben Hagerup, Peter Bro Miltersen, and Rasmus Pagh.
Deterministic dictionaries.
Journal of Algorithms, 41(1):69–85, 2001.

[3] Peter van Emde Boas.
Preserving order in a forest in less than logarithmic time.
In *16th Annual Symposium on Foundations of Computer Science (sfcs 1975)*,
pages 75–84. IEEE, 1975.

[4] Dan E Willard.
Log-logarithmic worst-case range queries are possible in space $\theta(n)$.
Information Processing Letters, 17(2):81–84, 1983.

- 1 Static dictionaries
- 2 Integer data structures
- 3 Graph data structures
- 4 Dynamization and persistence
- 5 Cache-oblivious structures
- 6 Succinct data structures
- 7 Computation in stream model
- 8 Literatura