

Deterministic Dictionaries

G.851

Lecture 15

André Schulz
notes

[Hagerup, Mitzenmacher, Pagh '91]

Goal : build a ^{static} dictionary for w -bit keys
in $O(u \log u)$ time on a
word-RAM using $O(u)$ space
deterministically!

History of deterministic dictionaries

- [Tarjan, Yao '79] $O(u)$ construction if w is $O(\log u)$
 u, w polynomially related
 $f(u)$ Universe: $\{0, \dots, 2^w - 1\}$
keys fit in a single word
(for $u \leq 2^w$) Table can be stored as bitvector (no many entries)
- [Fredman, Kolmós, Szemerédi '84]: deterministic FKS
needs $O(u^3 w)$
- [Raman '96]: Refinement of FKS idea: $O(u^2 w)$
- Remark: we can build $O(u)$ table & ds
in constant query time if
 $w \in O(u^{\Omega(1)})$ [Fusion Trees]
few entries
- [Alon, Naor '96]: $O(u w \log^4 u)$ but queries need $\Theta(\frac{w}{\log u})$

today $O(u \log u)$ construction } deterministic
 $O(d)$ look-up
 $O(u)$ space

Abstract ~~construction~~ construction

- ① Universe reduction ($n = u^{O(1)}$)
(Error correcting codes)
- ② Further universe reduction ($n = u^2$)
(Trees)
- ③ Solution for $n = u^2$
(Displacement method)

1. Universe reduction to $n = u^{O(1)}$

We want $\phi: \{0, 1\}^u \rightarrow \{0, 1\}^r$, such that

- $r = O(\log u)$
- ϕ is 1-1 on S

let $\psi: \{0, 1\}^u \rightarrow \{0, 1\}^{4u}$ be an error correcting code with relative minimum distance $\delta > 0$.

This means: $\forall x, y (x \neq y) : \cancel{x \& y \neq}$
 $\psi(x) \& \psi(y)$ differ at at ^{least} ~~at~~
 $4\delta w$ bits

Remark: $\delta \leq \frac{1}{2}$ can be picked for $w \geq 2$

Lemma There exist a set $D \subseteq [4w]$
for $\psi(x)$
of distinguishable bit positions (i.e.
 $\psi(x) \& \psi(y)$, $x \neq y$, differ at one of
the bits of D) with $|D| \leq 2 \log w / \log \frac{1}{1-\delta}$

Proof: We construct D incrementally:

$$D_0 \subseteq D_1 \subseteq D_2 \dots D \subseteq D_k \subseteq \{1, \dots, 4w\}$$

For D_j we have a natural partition

$$\text{of } S : \left\{ \cancel{C(D, v \in \{0,1\}^{|D|})} \right. \\ \left. \{ C(D, v), v \in \{0,1\}^{|D|} \} \right\} \text{ with}$$

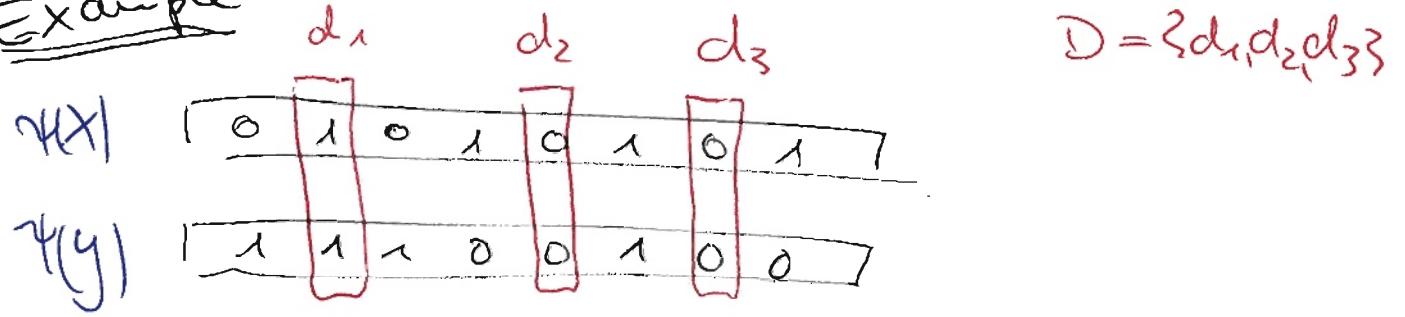
$$C(D, v) = \{ w \mid w \in S, \text{ the bits of } w \text{ at} \\ \text{the positions of } D \text{ match } v \}$$

↑
Cluster

$$B(D) = \sum_{K \in \{0,1\}^{100}} \binom{C(D, K)}{2} = \text{Collisions}$$

$\text{coll}(D_i) = \text{pairs not separated by } D_i$

Example



$\Rightarrow \{x, y\} \in C(\{d_1, d_2, d_3\}, 100)$

Fact: $\gamma(x)$ & $\gamma(y)$ can ~~be~~ water on at

most $2w(1-\delta)$ positions

For at most $2w(1-\delta)$ new bits x & y stay in the same

\Rightarrow



cluster

$$\sum_{\substack{\{x, y\} \in C(D_i, w) \\ x \neq y}} \sum_{d \in D} B(D_i \cup \{d\}) \leq \frac{2w(1-\delta)}{2w} B(D_i)$$

$$\Rightarrow \sum_{d=1}^{4w} \frac{B(D_i \cup \{d\})}{(1-\delta)4w} \leq \frac{1}{2} B(D_i) (1-\delta)$$

(merging principle)

$$\exists d : B(D_i \cup \{d\}) \leq \frac{B(D_i)}{2} (1-\delta)$$

Notice: it might be that $D_i \cup \{d\} = D_i$

We use $D_{i+1} = D_i \cup \{d\}$

Since $D_0 = \emptyset \Rightarrow B(D_0) \leq \binom{u}{2} \leq u^2$

we reach $B(D_k) < 1$ for

$$k = \lfloor 2 \log_{\frac{1}{2}} u \rfloor$$

$$\Rightarrow |D| < \cancel{O(u^2)} O(\log u) \quad \square$$

Lemma ~~Is~~ About position d :

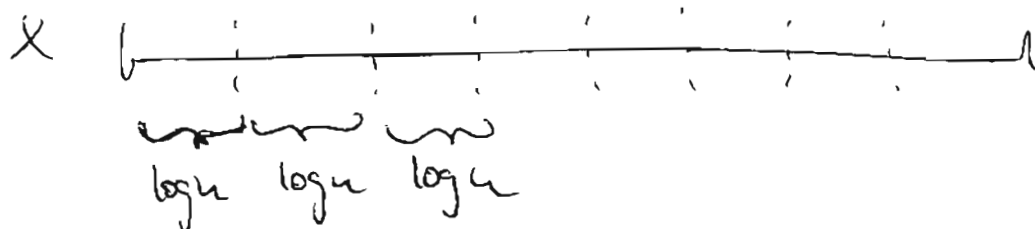
$$B(D \cup \{d\}) \leq (1 - \delta^{\frac{1}{2}}) B(D)$$

can be found in $O(u)$ time & space
deterministically

without proof (BIT-Tricks)

2. Universe reduction to $u = \lfloor n^2 \rfloor$

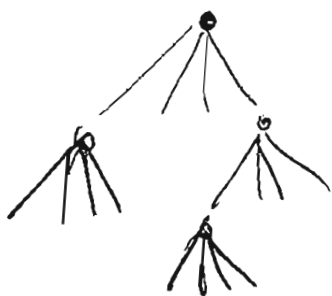
split $x \in u$ into $\log u$ chunks



$\Theta(1)$ many
chunks

"# of $\log u$ Blocks" = $n \Rightarrow$
possible

Store $x \in S$ in a tree with $|\Sigma| = u$



• Fan-out $\leq u$

• nodes: u

↳ Every node gets an ID

To find a node we build a table that can answer queries of the form

$(ID, a) \rightarrow$ which node can I reach

from node ID following $a \in \Sigma^*$

$O(u^2)$ tuples to store in hash-table

We also store in every node of the tree

if it is terminal

\Rightarrow We can find $x \in S$ ^{and observe if} $x \notin S$ using R queries

\Rightarrow Search time is $O(1)$ because the

height of the tree is $O(1)$

3. Solution for a quadratic universe

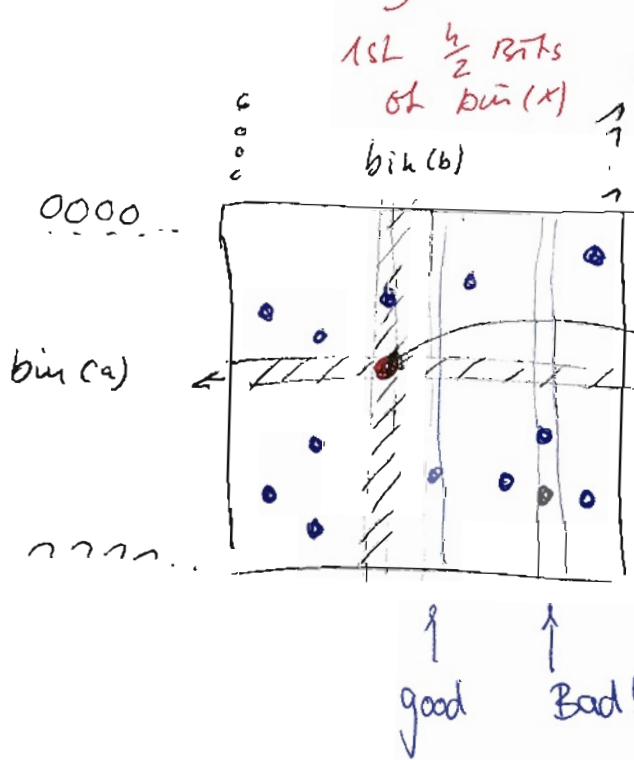
We search a perfect hash function

$$h: [u^2] \rightarrow [2^r], \quad r = \log u + o(1)$$

We first obtain a function

$$\bar{h}: [u^2] \rightarrow [2^r]^2$$

$$\bar{h}(x) := (\text{low}(x), \text{high}(x)), \quad \bar{h}(x) \text{ is 1-1 on } S$$



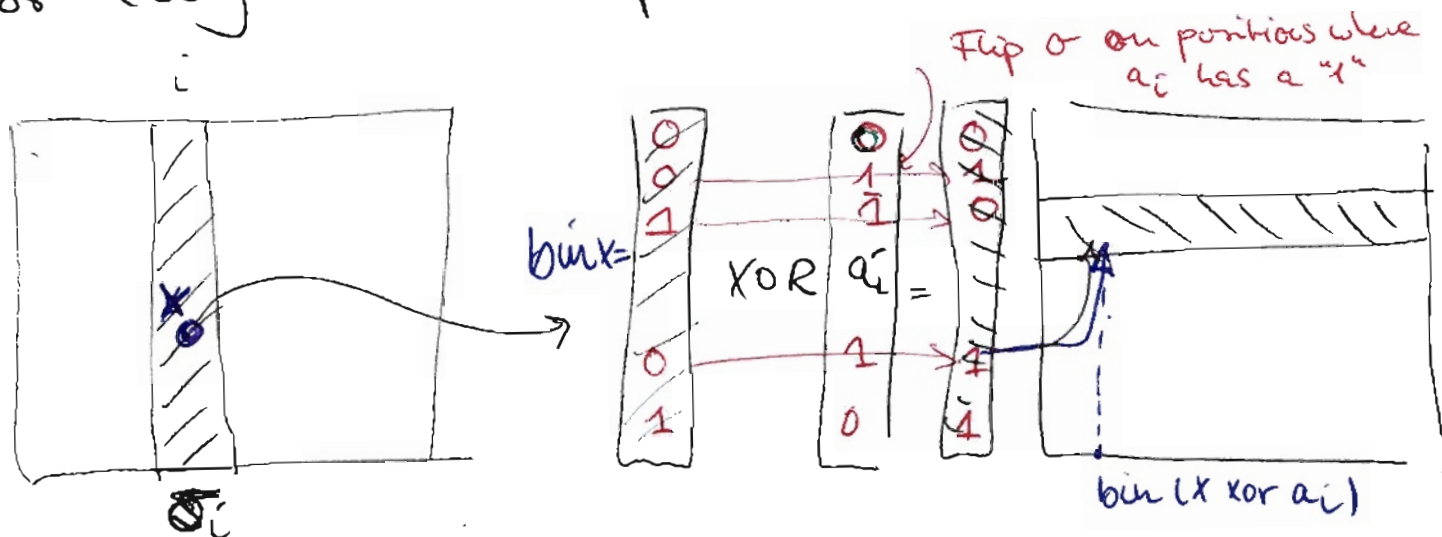
$$\bar{h}(x) = (a, b)$$

if every column would contain ~~one~~ only 1 entry of $\{\bar{h}(x) \mid x \in S\}$

we can use the column # as hash value

Idea: Scramble the table

For every column i pick some $a_i \in \{0,1\}^r$



→ compute $\sigma_i \text{ XOR } a_i$ and add the vector as row i in the new table

Formal: $(x, y) \in \text{Table} \rightarrow (y \text{ XOR } a_x, x) \in \text{NewTable}$

f = collisions in the columns of Table

f' = collisions in the columns of NewTable

Lemma There is a set of $\{a_i\}_{i=1}^r$ such that $f' \leq \min \{u, \lfloor \frac{f}{2^{r-3}} \rfloor u\}$

Proof: (1) Reorder columns, such that

$$|S_1| \geq |S_2| \geq \dots \geq |S_r|$$

↑ S_i = ^{rows} elements in the i th column of Table

-8- $S_i = \{x \mid (i, x) \in \bar{u}(S)\}$

• Notice: $q = \sum_i \binom{|S_i|}{2}$

• (2) Pick $a_1, a_2, a_3, \dots, a_{2^r}$ in increasing order

Assume we ~~want~~ want to find a_i

→ pick a_i uniformly at random from $\mathbb{Z}_{0,1}^{2^r}$

↳ let (m_j) ^{Def.} be the collisions in NewTable at column j so far

~~new collisions~~

new collisions by picking $a_i = \sum_{y \in S_i} m_y \cdot \overbrace{y}^{\text{def.}} \cdot a_i$

$$\begin{aligned} \rightarrow E\left(\sum_{\substack{y \in S_i \\ y}} m_y \cdot y \cdot a_i\right) &= \sum_{y \in S_i} E(m_y \cdot y \cdot a_i) \\ &= \sum_y \sum_{\substack{u \in \mathbb{Z}_{0,1}^{2^r} \\ u}} \underbrace{\mathbb{P}[(y \cdot x \cdot a_i) = u]}_{\frac{1}{2^r}} \cdot m_y \cdot u \\ &= \sum_{y \in S_i} \frac{1}{2^r} \sum_{u \in \mathbb{Z}_{0,1}^{2^r}} m_y \cdot u \\ &= \frac{|S_i| \sum_{j \in C} |S_j|}{2^r} = \sum_{j < i} |S_j| \quad (*) \end{aligned}$$

By Markov's inequality:

$$\mathbb{P}[\text{new collisions} \leq 2 \cdot (*)] \geq \frac{1}{2}$$

↳ After 2 tries we found a a_i vector:

$$\# \text{ new collisions} \leq \lfloor 2^{\lceil r \rceil} |S_i| \sum_{j < i} |S_j| \rfloor = M_i$$

↑ floor, because # is $\in \mathbb{Z}$

$$\Rightarrow q' \leq \sum_i \left\lfloor |S_i| \sum_{j < i} |S_j| \cdot 2^{-r+1} \right\rfloor$$

Consequence for $r = \log u + 4$

$$(A) \quad q' \leq \sum_i |S_i| n \cdot 2^{-r+1} \leq n 2^{-r+1} \leq n$$

$$(B) \quad q' \leq \sum_i M_i$$

$$(B.1) \quad |S_i| = 1 \Rightarrow M_i = \lfloor 2 \cdot \sum_{j < i} \frac{|S_j|}{2^r} \rfloor \leq \lfloor 2 \cdot \frac{n}{2^r} \rfloor = 0$$

$$(B.2) \quad |S_i| > 1 \Rightarrow \forall j < i: |S_j|^2 \leq 4 \left(\frac{|S_j|}{2} \right)$$

$$\Rightarrow M_i \leq \lfloor 2 \cdot \sum_{j < i} \frac{|S_j|^2}{2^r} \rfloor \leq \lfloor 2^{-r+3} \sum_{j < i} |S_j| \rfloor$$

$$\Rightarrow q' \leq \sum_i \lfloor 2^{-r+3} q' \rfloor \leq n \cdot \lfloor 2^{-r+3} q' \rfloor$$

□

⇒ 1st ~~shuffling~~ shuffle

lemma

↳ $q' \leq n$

2nd ~~shuffling~~ shuffle

lemma

↳ $q^n \leq \left\lfloor \frac{q}{2^{r+3}} \right\rfloor \cdot n = \left\lfloor \frac{n}{2^{\log n + 2}} \right\rfloor \cdot n = \underline{\underline{0}}$

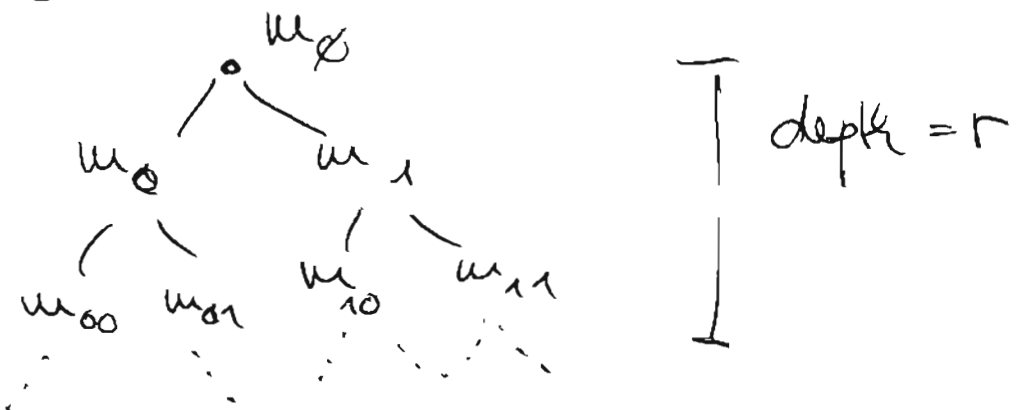
Derandomization of the selection of a_i

• Extension of notation

$\forall z \in \{0, 1\}^k$ \uparrow $m_z = \sum_{v \in \{0, 1\}^{r-k}} m_{z \circ v}$

• We want to keep track of all m_x during the ~~search~~ search for a_i

↳ Store current m_x values inside a binary tree



→ ~~Adding~~ New a_i value → update tree
↳ insert all $y \in S_i$ in the tree by
computing $a_i \cdot y$ and increase
all w_x 's on the associated
root → leaf path

$O(|S_i| \cdot r)$ for every a_i

⇒ $O(r \cdot n) = O(n \log n)$ total

Derandomize by the technique of

CONDITIONAL EXPECTATIONS

IDEA • We know that $E(q')$ is "good" ($\leq \epsilon$)
for some a_i we have to pick

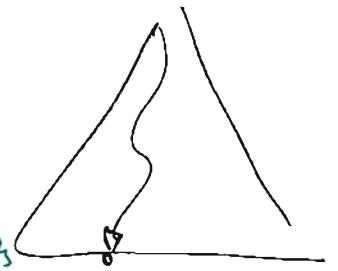
→ ~~Either~~

$E(q' / a_i \text{ starts with } z)$ is also
good for either $z=0$ or $z=1$

→ We use the tree to determine
the value of z and determine
 a_i bit by bit

$$E(q') = \sum_{y \in S_i} \sum_{n \in \{0,1\}^r} w_n \frac{1}{2^r} = \sum_{y \in S_i} \frac{1}{2^r} w_y$$

maybe use $S_i^{(y)} = \{x \in S_i \mid x \neq z \dots\}$
notation



$$E(q' \mid a_i \text{ starts with } 0) = \sum_{\substack{y \in S_i \\ y=1\dots}} \sum_{n \in \{0,1\}^{r-1}} \frac{1}{2^{r-1}} w_{yn} + \sum_{\substack{y \in S_i \\ y=0\dots}} \sum_{n \in \{0,1\}^{r-1}} \frac{1}{2^{r-1}} w_{0n}$$

$$= \frac{1}{2^{r-1}} \left(|\{y \in S_i \mid y=1\dots\}| \cdot w_1 + |\{y \in S_i \mid y=0\dots\}| \cdot w_0 \right)$$

$$E(q' \mid a_i \text{ starts with } 1) = \frac{1}{2^{r-1}} \left(|\{y \in S_i \mid y=0\dots\}| w_0 + |\{y \in S_i \mid y=1\dots\}| w_1 \right)$$

Check which one is ^{smaller} larger and use the corresponding bit in a_i

↳ iterate

y in the end

↳ no randomness here

$$E(q' \mid a_i = \underbrace{01001\dots 01}_{\text{something like that}}) \neq (*)$$

$$\Rightarrow q' \text{ for } a_i = 01\dots \leq (*)$$

Idea dynamic perfect Hashing

[Dietzfelbinger et al 1994]

After

- Every $w' = c \cdot n$ insert/delete opn rebuild
 - Delete: mark a node as deleted if it is stored
 - insert: Lookup node: if marked as deleted, unmark
if not there:
 - Add in table
 - If collision, rehash 2nd level table
 - If too many elements in the 2nd level table
 - ↳ rehash everything
- ⇒ $O(w)$ space, $O(1)$ per opn in expectation & amortized