

CS 240 – Data Structures and Data Management

Module 7E: Hashing - Enriched

T. Biedl E. Kondratovsky M. Petrick O. Veksler

Based on lecture notes by many previous cs240 instructors

David R. Cheriton School of Computer Science, University of Waterloo

Winter 2022

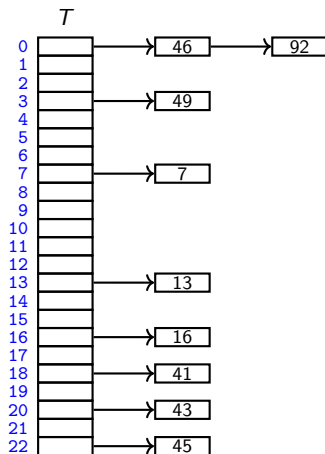
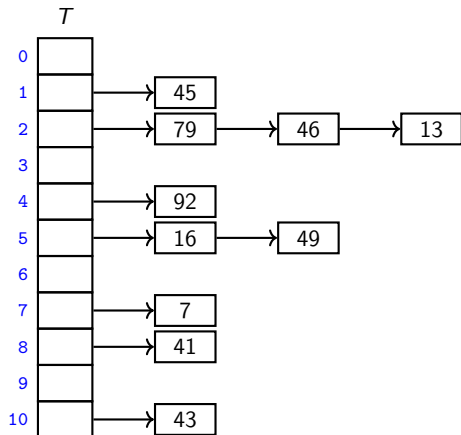
Outline

- Rehashing
- Multiplication method
- Randomly chosen hash-functions

Outline

- Rehashing
- Multiplication method
- Randomly chosen hash-functions

Rehashing



Outline

- Rehashing
- **Multiplication method**
- Randomly chosen hash-functions

Multiplication method

- Pick some number $A \in (0, 1)$ (preferably an irrational)

$$h(k) = \left\lfloor M \cdot \left(\underbrace{A \cdot k}_{\text{multiply}} - \underbrace{\lfloor A \cdot k \rfloor}_{\text{integral part}} \right) \right\rfloor$$

$\underbrace{\hspace{15em}}_{\text{fractional part, in } [0, 1]}$
 $\underbrace{\hspace{25em}}_{\text{integer in } [0, M]}$

- Example:

$$A = 0.a_1 a_2 a_3 \dots$$

$$k = b_1 b_2 \dots b_6$$

(both in base 2)

$A \cdot k$	=	0 0 0 .. 0 0	(leading bits)		(bits of fractional part)
$+a_1 \cdot$		0 $b_1 b_2 b_3$.. b_5		b_6	
$+a_2 \cdot$		0 0 $b_1 b_2 b_3$..		$b_5 b_6$	
$+a_3 \cdot$		0 0 0 $b_1 b_2 b_3$.. $b_5 b_6$	
\vdots		\vdots		\vdots	
$+a_5 \cdot$		0 0 0 0 0 b_1		$b_2 b_3$..	$b_5 b_6$
$+a_6 \cdot$		0 0 0 0 0 0		$b_1 b_2 b_3$.. $b_5 b_6$
$+a_7 \cdot$		0 0 0 0 0 0		0 $b_1 b_2$	b_3 .. $b_5 b_6$
$+a_8 \cdot$		0 0 0 0 0 0		0 0 b_1	$b_2 b_3$.. $b_5 b_6$
\vdots		\vdots		$\leftarrow h(k) \rightarrow$	\vdots

- Should use at least $\log |U| + \log |M|$ bits of A .

Outline

- Rehashing
- Multiplication method
- Randomly chosen hash-functions

Randomly chosen hash-functions

- There are too many possible hash-functions \Rightarrow we cannot compute hash-value quickly if we choose randomly among them.
- Idea: fix a family \mathcal{H} of hash-functions that are easy to compute.
- But what should criteria be for \mathcal{H} ?
- Uniform hash-values, i.e., $P(h(k) = i) = \frac{1}{M}$, is *not* enough.

\mathcal{H}_1	keys		
	x	y	z
h_0	0	0	0
h_1	1	1	1

- $M = 2$ in this example.
- $P(h(k) = i) = \frac{1}{2}$ for $i = 0, 1$ and $k = x, y, z$
- But these hash-functions are terrible!

- Goal: Small probability of collisions (**universal hashing**):

$$P(h(k) = h(k')) = \frac{1}{M} \quad \text{for any two keys } k \neq k'$$

- This is enough for hashing-analysis for chaining to hold.

Carter-Wegman hash-function

$$h_{a,b}(k) = \left(\underbrace{a \cdot k + b \bmod p}_{f_{a,b}(k)} \right) \bmod M$$

(where $k \in \mathbb{Z}_p$, p prime, $a, b \in \mathbb{Z}_p$ chosen randomly, $a \neq 0$)

Example: ($p = 5$, $M = 2$):

	keys				
	0	1	2	3	4
$f_{1,0}$	0	1	2	3	4
$f_{2,0}$	0	2	4	1	3
$f_{1,2}$	2	3	4	0	1
$f_{2,1}$	1	3	0	2	4
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

Carter-Wegman hash-function

$$h_{a,b}(k) = \left(\underbrace{a \cdot k + b \bmod p}_{f_{a,b}(k)} \right) \bmod M$$

(where $k \in \mathbb{Z}_p$, p prime, $a, b \in \mathbb{Z}_p$ chosen randomly, $a \neq 0$)

Example: ($p = 5$, $M = 2$):

	keys				
	0	1	2	3	4
$f_{1,0}$	0	1	2	3	4
$f_{2,0}$	0	2	4	1	3
$f_{1,2}$	2	3	4	0	1
$f_{2,1}$	1	3	0	2	4
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

	keys				
	0	1	2	3	4
$h_{1,0}$	0	1	0	1	0
$h_{2,0}$	0	0	0	1	1
$h_{1,2}$	0	1	0	0	1
$h_{2,1}$	1	1	0	0	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

Carter-Wegman hash-function

$$h_{a,b}(k) = \left(\underbrace{a \cdot k + b \bmod p}_{f_{a,b}(k)} \right) \bmod M$$

(where $k \in \mathbb{Z}_p$, p prime, $a, b \in \mathbb{Z}_p$ chosen randomly, $a \neq 0$)

Example: ($p = 5$, $M = 2$):

	keys				
	0	1	2	3	4
$f_{1,0}$	0	1	2	3	4
$f_{2,0}$	0	2	4	1	3
$f_{1,2}$	2	3	4	0	1
$f_{2,1}$	1	3	0	2	4
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

	keys				
	0	1	2	3	4
$h_{1,0}$	0	1	0	1	0
$h_{2,0}$	0	0	0	1	1
$h_{1,2}$	0	1	0	0	1
$h_{2,1}$	1	1	0	0	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

Claim: $f_{a,b}$ is a permutation of \mathbb{Z}_p .