

# 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

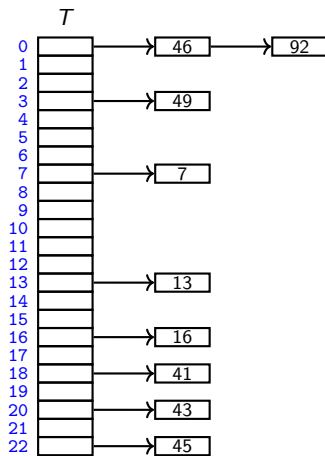
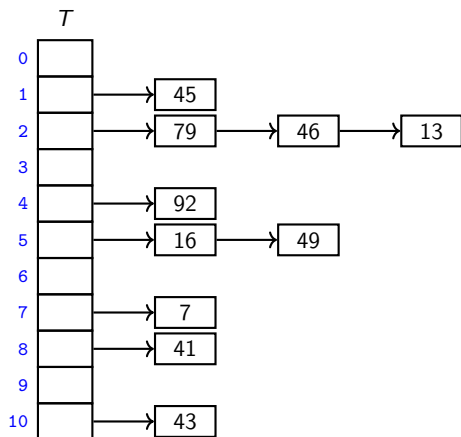
# Outline

- 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  $A \in (0, 1)$  (preferably an irrational, e.g.  $A = \frac{\sqrt{5}-1}{2} = \phi$ )

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

fractional part, in  $[0, 1)$

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)

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

	$A \cdot k$	$=$	(leading bits) 0 0 0 .. 0 0	(bits of fractional part)	
	$+ a_1 \cdot$		0 b <sub>1</sub> b <sub>2</sub> b <sub>3</sub> .. b <sub>5</sub>	b <sub>6</sub>	
	$+ a_2 \cdot$		0 0 b <sub>1</sub> b <sub>2</sub> b <sub>3</sub> ..	b <sub>5</sub> b <sub>6</sub>	
	$+ a_3 \cdot$		0 0 0 b <sub>1</sub> b <sub>2</sub> b <sub>3</sub>	.. b <sub>5</sub> b <sub>6</sub>	
	$\vdots$			..	
	$+ a_5 \cdot$		0 0 0 0 0 b <sub>1</sub>	b <sub>2</sub> b <sub>3</sub> ..	b <sub>5</sub> b <sub>6</sub>
	$+ a_6 \cdot$		0 0 0 0 0 0	b <sub>1</sub> b <sub>2</sub> b <sub>3</sub>	.. b <sub>5</sub> b <sub>6</sub>
	$+ a_7 \cdot$		0 0 0 0 0 0	0 b <sub>1</sub> b <sub>2</sub>	b <sub>3</sub> .. b <sub>5</sub> b <sub>6</sub>
	$+ a_8 \cdot$		0 0 0 0 0 0	0 0 b <sub>1</sub>	b <sub>2</sub> b <sub>3</sub> .. b <sub>5</sub> b <sub>6</sub>
	$\vdots$			$\leftarrow h(k) \rightarrow$	..

# Outline

- Rehashing
- Multiplication method
- Randomly chosen hash-functions

# Randomly chosen hash-functions

- There are  $|U|^M$  many possible hash-functions
- Ideally we would choose randomly among all of them.
- But then we cannot compute hash-value quickly!
- **Idea:** Fix a family  $\mathcal{H}$  of hash-functions that are easy to compute. Then choose uniformly among them.
- **Example:**

- ▶  $U = \mathbb{Z}_5, M = 2$
- ▶  $h_b(k) = ((k + b) \bmod 5) \bmod 2$
- ▶  $\mathcal{H} = \{h_b : b \in \mathbb{Z}_5\}$
- ▶ Choose  $b \in \mathbb{Z}_5$  randomly to get hash-function

$\mathcal{H}$	keys				
	0	1	2	3	4
$h_0$	0	1	0	1	0
$h_1$	1	0	1	0	0
$h_2$	0	1	0	0	1
$h_3$	1	0	0	1	0
$h_4$	0	0	1	0	1

- But how do we measure whether these are “good”?



# Universal hash-functions

- For analysis, we needed *uniform hash-values*:

$$P(h(k) = i) = \frac{1}{M}$$

- But this is *not* good enough.

	keys				
$\mathcal{H}$	0	1	2	4	5
$h_0$	0	0	0	0	0
$h_1$	1	1	1	1	1

- $P(h(k) = i) = \frac{1}{2}$  for  $i = 0, 1$  and any  $k$
- But these hash-functions are terrible!
- Problem: hash-values not independent

- Also want: 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 analyzing hashing with chaining as before.

# 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,  $a \neq 0$ ),  $M < p$

**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  $\mathbb{Z}_p$ .