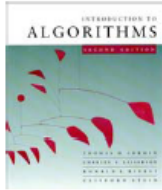# Universal Hashing

# A weakness of hashing
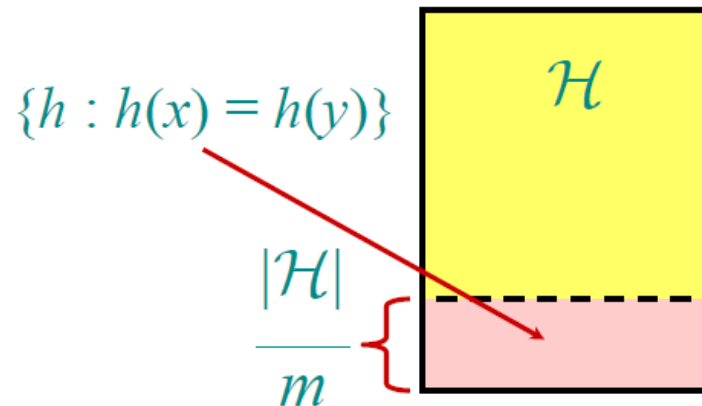
- ## Problem

  - For any hash function h, a set of keys exist that can cause the average access time of a hash table to skyrocket.

    - An adversary can pick all the keys which all map to the same bucket

- ## Idea

  - Choose the hash function at random, independent of the keys
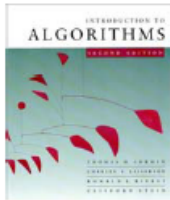
# Universal hashing

**Definition.** Let $U$ be a universe of keys, and let $\mathcal{H}$ be a finite collection of hash functions, each mapping $U$ to $\{0, 1, \ldots, m-1\}$. We say $\mathcal{H}$ is *universal* if for all $x, y \in U$, where $x \neq y$, we have $|\{h \in \mathcal{H} : h(x) = h(y)\}| = |\mathcal{H}|/m$.

That is, the chance of a collision between $x$ and $y$ is $1/m$ if we choose $h$ randomly from $\mathcal{H}$.

$\{h : h(x) = h(y)\}$

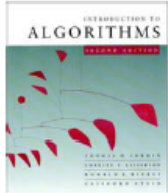$\mathcal{H}$

$\dfrac{|\mathcal{H}|}{m}$

# Universality is good

**Theorem.** Let $h$ be a hash function chosen (uniformly) at random from a universal set $\mathcal{H}$ of hash functions. Suppose $h$ is used to hash $n$ arbitrary keys into the $m$ slots of a table $T$. Then, for a given key $x$, we have
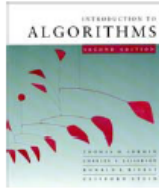
$$E[\#\text{collisions with } x] < n/m.$$

# Proof of theorem

*Proof.* Let $C_x$ be the random variable denoting the total number of collisions of keys in $T$ with $x$, and let

$$c_{xy} = \begin{cases} 1 & \text{if } h(x) = h(y), \\ 0 & \text{otherwise.} \end{cases}$$

Note: $E[c_{xy}] = 1/m$ and $C_x = \displaystyle\sum_{y \in T-\{x\}} c_{xy}$ .
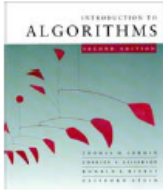
# **Proof (continued)**

$$E[C_x] = E\left[\sum_{y \in T - \{x\}} c_{xy}\right]$$

- Take expectation of both sides.

$$= \sum_{y \in T - \{x\}} E[c_{xy}]$$

- Linearity of expectation.

$$= \sum_{y \in T - \{x\}} 1/m$$

- $E[c_{xy}] = 1/m$.

$$= \frac{n-1}{m} \cdot \;\blacksquare$$

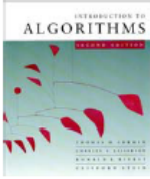- Algebra.

# Constructing a set of universal hash functions

Let $m$ be prime. Decompose key $k$ into $r + 1$ digits, each with value in the set $\{0, 1, \ldots, m{-}1\}$. That is, let $k = \langle k_0, k_1, \ldots, k_r \rangle$, where $0 \le k_i < m$.

**Randomized strategy:**

Pick $a = \langle a_0, a_1, \ldots, a_r \rangle$ where each $a_i$ is chosen randomly from $\{0, 1, \ldots, m{-}1\}$.

Define $h_a(k) = \displaystyle\sum_{i=0}^{r} a_i k_i \bmod m$.   *Dot product, modulo $m$*

How big is $\mathcal{H} = \{h_a\}$?   $|\mathcal{H}| = m^{r+1}$.  ← **REMEMBER THIS!**

# Universality of dot-product hash functions
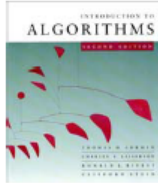
**Theorem.** The set $\mathcal{H} = \{h_a\}$ is universal.

*Proof.* Suppose that $x = \langle x_0, x_1, \ldots, x_r \rangle$ and $y = \langle y_0, y_1, \ldots, y_r \rangle$ be distinct keys. Thus, they differ in at least one digit position, wlog position $0$. For how many $h_a \in \mathcal{H}$ do $x$ and $y$ collide?

We must have $h_a(x) = h_a(y)$, which implies that

$$\sum_{i=0}^{r} a_i x_i \equiv \sum_{i=0}^{r} a_i y_i \quad (\mathrm{mod}\, m) .$$

# Proof (continued)

Equivalently, we have

$$\sum_{i=0}^{r} a_i(x_i - y_i) \equiv 0 \quad (\bmod\, m)$$

or

$$a_0(x_0 - y_0) + \sum_{i=1}^{r} a_i(x_i - y_i) \equiv 0 \quad (\bmod\, m)\,,$$

which implies that

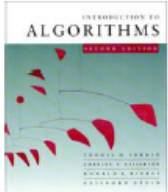$$a_0(x_0 - y_0) \equiv -\sum_{i=1}^{r} a_i(x_i - y_i) \quad (\bmod\, m)\,.$$

# Fact from number theory

**Theorem.** Let $m$ be prime. For any $z \in \mathbb{Z}_m$ such that $z \neq 0$, there exists a unique $z^{-1} \in \mathbb{Z}_m$ such that

$$z \cdot z^{-1} \equiv 1 \pmod{m}.$$

**Example:** $m = 7$.

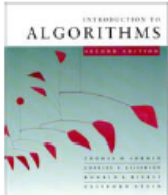| $z$ | 1 | 2 | 3 | 4 | 5 | 6 |
|------|---|---|---|---|---|---|
| $z^{-1}$ | 1 | 4 | 5 | 2 | 3 | 6 |

# Back to the proof

We have

$$a_0(x_0 - y_0) \equiv -\sum_{i=1}^{r} a_i(x_i - y_i) \quad (\mathrm{mod}\, m),$$

and since $x_0 \neq y_0$, an inverse $(x_0 - y_0)^{-1}$ must exist, which implies that

$$a_0 \equiv \left( -\sum_{i=1}^{r} a_i(x_i - y_i) \right) \cdot (x_0 - y_0)^{-1} \quad (\mathrm{mod}\, m).$$

Thus, for any choices of $a_1, a_2, \ldots, a_r$, exactly one choice of $a_0$ causes $x$ and $y$ to collide.

# Proof (completed)

**Q.** How many $h_a$'s cause $x$ and $y$ to collide?

**A.** There are $m$ choices for each of $a_1, a_2, \ldots, a_r$, but once these are chosen, exactly one choice for $a_0$ causes $x$ and $y$ to collide, namely

$$a_0 = \left( \left( -\sum_{i=1}^{r} a_i(x_i - y_i) \right) \cdot (x_0 - y_0)^{-1} \right) \bmod m.$$

Thus, the number of $h$'s that cause $x$ and $y$ to collide is $m^r \cdot 1 = m^{r} = |\mathcal{H}|/m$. ▨