# COL 351: Analysis and Design of Algorithms

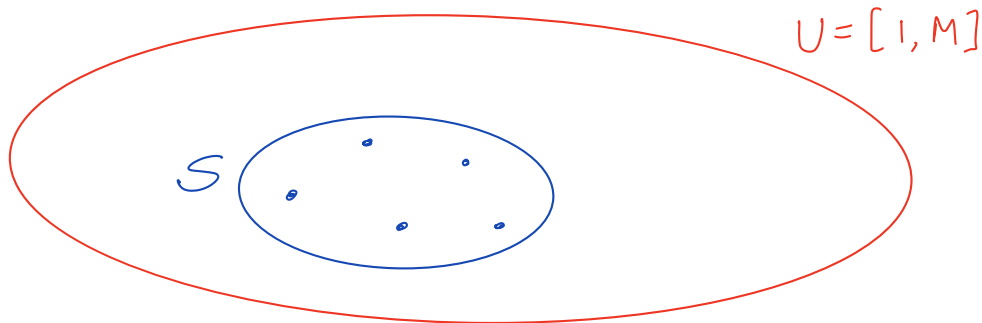## Lecture 20

# Set Membership

**Given:** A universe $U = [1, 2, \ldots, M]$, and a set $S \subsetneq [1,M]$ of size $n$.

**Goal:** Find a data-structure of $O(n = |S|)$ size that answers for any $x \in [1,M]$ query of form:

"Does $x \in S$ ?"

$$U = [1, M]$$

# Hash Function

$$H(z) = z \mod n$$

- Works well for a random $S$

- What if $S$ is not random?

| |
|---|
| 0 |
| |
| $i$ |
| |
| |
| $n-1$ |

Table $T$

$z_1 \rightarrow z_2 \rightarrow z_3$

Store in link-list the set

$\{x \in S \mid H(x) = i\}$

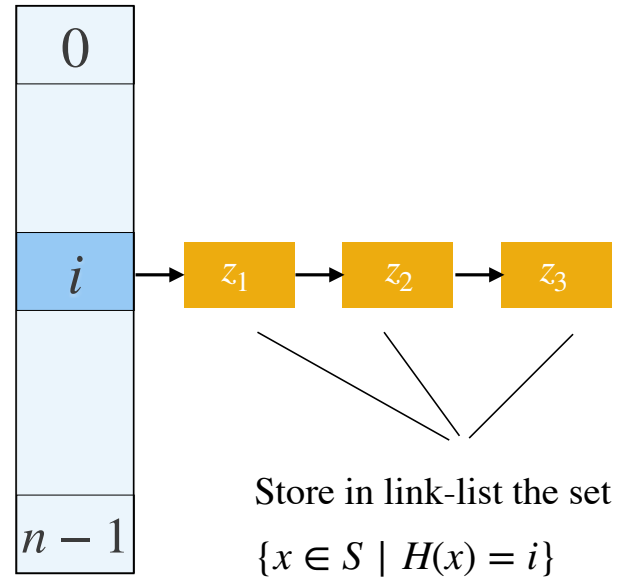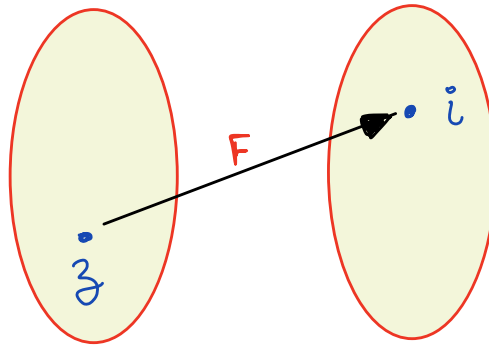**Claim:** No single hash function can work for all possible sets S

# Modular Arithmetic

$$F(z) = (r \cdot z) \mod p \qquad (\text{Here, } p \text{ is a prime}).$$

**Claim:** If $r \in [1, p-1]$ was random, then for any $z, i \in [1, p-1]$, we have
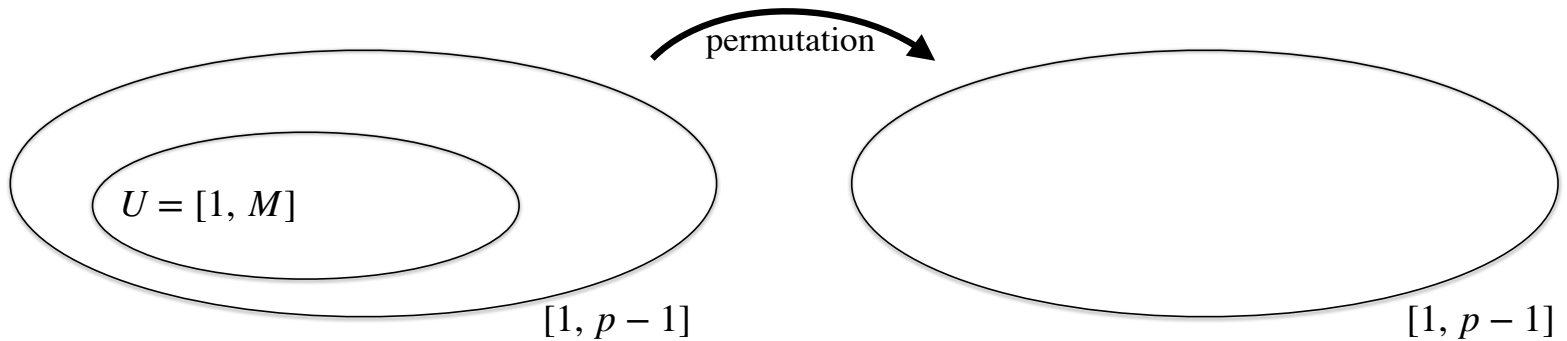
$$\text{Prob}(F(z) = i) = \frac{1}{p-1}.$$

# New Hash Function

- Universe $U = [1, M]$.

- $p =$ prime in range $[M+1, 2M]$, $\quad r =$ integer range $[1, p-1]$

> **Hash Function:**
>
> $H_r(z) \;=\; \big((r \cdot z) \mod p\big) \mod n$



permutation

$U = [1, M]$

$[1, p-1]$

$[1, p-1]$

# New Hash Function

**Hash Function:**

$$H_r(z) = \big((r \cdot z) \mod p\big) \mod n$$

*What is collision probability?*

| |
|---|
| 0 |
| |
| $i$ |
| |
| $n-1$ |

$z_1 \rightarrow z_2 \rightarrow z_3$

Table $T$
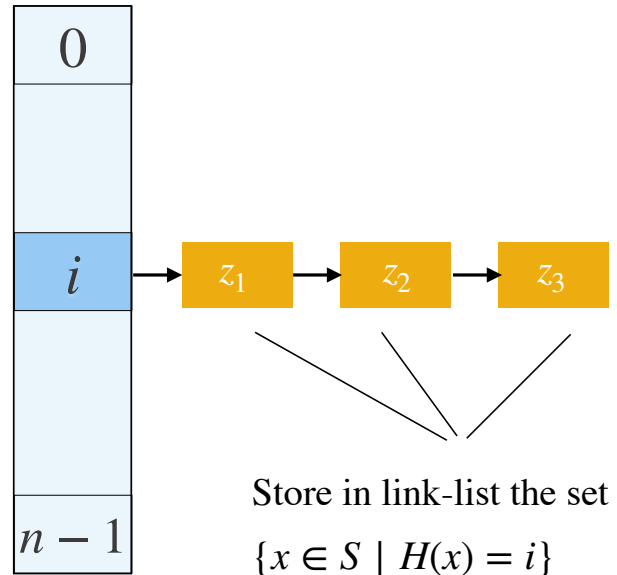
Store in link-list the set

$\{x \in S \mid H(x) = i\}$

**Hash Function:**

$$H_r(z) = \left((r \cdot z) \bmod p\right) \bmod n$$

*Question:* For any $x, y \in U$, what is **collision probability** if '$r$' is randomly chosen?

Solution:

$H_r(x) = H_r(y)$

$\Rightarrow (rx \bmod p) - (ry \bmod p)$ is multiple of $n$

$\Rightarrow (rx - ry \bmod p)$ is multiple of $n$    **or**    $(rx - ry \bmod p) - p$ is multiple of $n$

$\Rightarrow (rx - ry \bmod p)$ lies in $\{n, 2n, 3n, \ldots\}$   **or**   $\{p - n, p - 2n, p - 3n, \ldots\}$

$\Rightarrow (r(x - y) \bmod p)$ lies in $\{n, 2n, 3n, \ldots, p - 3n, p - 2n, p - n\}$

$$\text{Prob}\left(H_r(x) = H_r(y)\right) \leq \frac{1}{p-1} \cdot |\{n, 2n, 3n, \ldots, p - 3n, p - 2n, p - n\}| \approx \frac{2}{n}$$

# Expected Time to search an element

*Question*: For any $x \in U$, what is expected time to verify membership of $x$ in set $S$?

Solution:

The time to search $x$ is sum of

    (i) Time to compute $H_r(x)$, and

    (ii) Number of elements in $S$ mapped to $H_r(x)$.

Expected Time:

$$= 1 + \sum_{y \in S \setminus \{x\}} \text{Prob}\big(H_r(y) = H_r(x)\big) \leq 1 + (n-1) \cdot \frac{2}{n} = O(1)$$

# Total number of Collisions

*Question*: What is expected number of total collisions?

Solution:

Expected total number of collisions are

$$= \sum_{\substack{x,y \in S \\ x \neq y}} \mathrm{Prob}\big(H_r(y) = H_r(x)\big) \leq \frac{n(n-1)}{2} \cdot \frac{2}{n} \leq n$$

# Balls – Bins Exercise

n balls    ◯  ◯  ◯  - - -  ◯

n bins    ⊔  ⊔  ⊔  - - -  ⊔

- each ball goes into one of the randomly selected Bin

- Expected no of balls in Bin i $= \frac{1}{n} \cdot n = 1$

- <u>Fact</u>

$$Exp\left( \max_{i=1}^{n} (\# \text{ of Balls in Bin } i)\right) = \Theta\left(\frac{\log n}{\log \log n}\right)$$

$sol^n$ to $R^R = n$

# Maximum Time to search an element

Worst – Case time

**Question**: What is expected value of $\max\limits_{i\in[0,\,n-1]} |T[i]|$?

Answer:

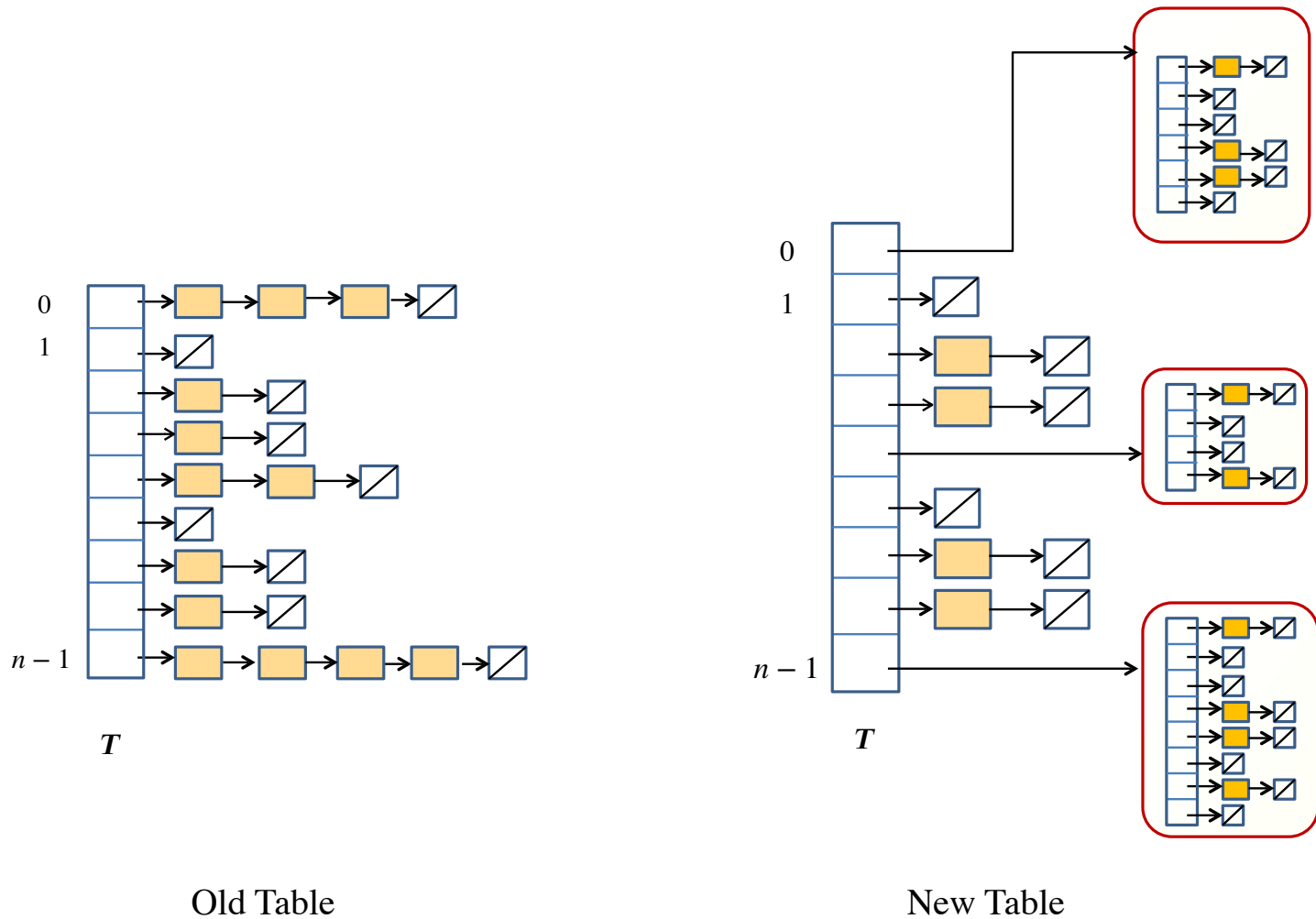$$\text{It will be } \Theta\left(\frac{\log n}{\log\log n}\right).$$

*Is hashing any better than AVL trees?*

Remark:

The proof of the fact $\text{Exp}\left(\max\limits_{0\le i\le n-1} |T[i]|\right) = \Theta\left(\frac{\log n}{\log\log n}\right)$

is not part of Syllabus.

# Two-Level Hash Table



Old Table

New Table

# Two-Level Hash Table

**Outer Hash Function:**

$$H_r(z) = \big((r \cdot z) \bmod p\big) \bmod n$$

**Inner Hash Function:**

$$z \mapsto \big((r_0 \cdot z) \bmod p\big) \bmod n_i^2$$

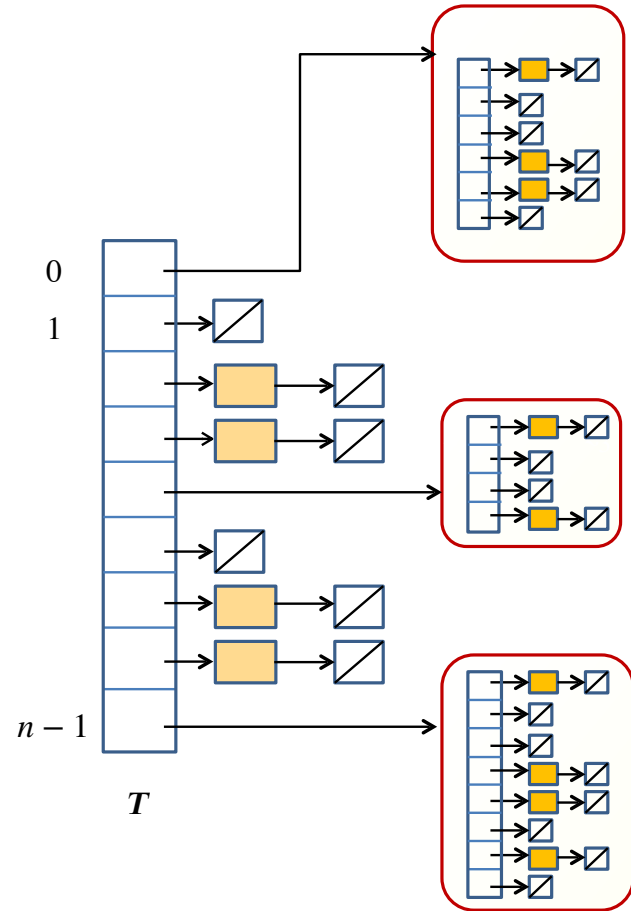where, $n_i$ = size of $T[i]$



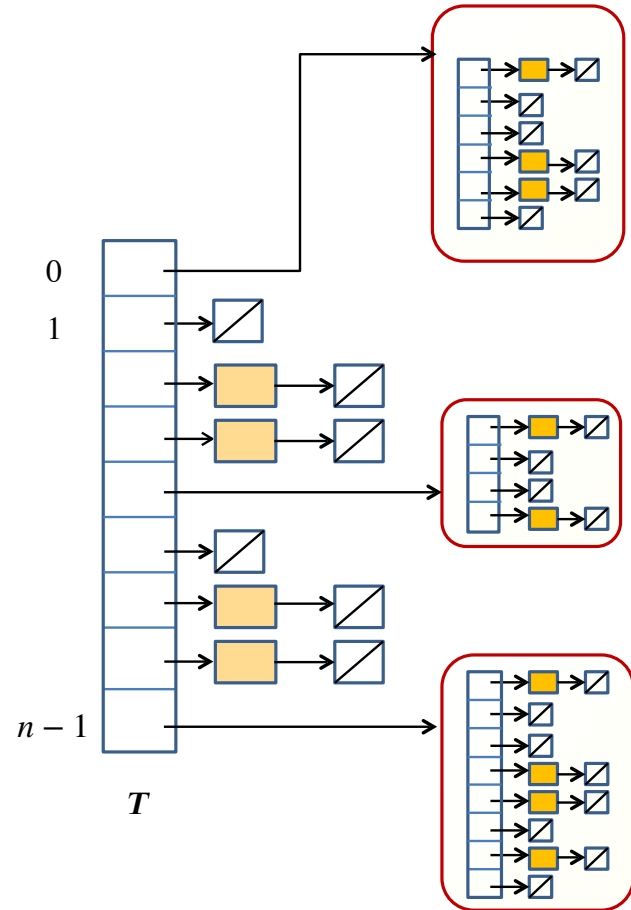New Table

# Two-Level Hash Table

**Outer Hash Function:**

$$H_r(z) = ((r \cdot z) \mod p) \mod n$$

**Inner Hash Function:**

$$z \mapsto ((r_0 \cdot z) \mod p) \mod n_i^2$$

where, $n_i$ = size of $T[i]$



$T$

New Table

- What is expected total size?
- What is expected number of total collisions?