# Optimal Dynamic Subset Sampling: Theory and Applications

Lu Yi, Hanzhi Wang, Zhewei Wei    Contact: zhewei@ruc.edu.cn

RENMIN UNIVERSITY OF CHINA

KDD2023 · AUGUST 6–10 · LONG BEACH, CA

## Overview

**Subset Sampling Problem**

➤ Given a set of $n$ distinct events $S = \{x_1, \ldots, x_n\}$, in which each event $x_i$ has an associated probability $p(x_i)$, a query for the subset sampling problem returns a subset $T \subseteq S$, such that every $x_i$ is independently included in $T$ with probability $p(x_i)$.

domain set $S$

$x_1$ $x_2$ $x_3$ $x_4$ ... $x_n$

$p(x_1)$ $p(x_2)$ $p(x_3)$ $p(x_4)$ $p(x_n)$

$\sum_i p(x_i) = \mu$

query →

sample result $T \subseteq S$

$x_1$ $x_4$

Each $x_i$ is included in $T$ independently with probability $p(x_i)$

**Dynamic Subset Sampling Problem**

➤ Insert an event
➤ Delete an event
➤ Modify the probability of an event

Influence spreading under the IC (Independent Cascade) model

⭐⭐ **Contributions**

✓ **Optimal** query time: $O(1 + \mu)$
✓ **Optimal** update time: $O(1)$
✓ Great **experimental performance**
✓ Empirical study on **Influence Maximization**

**Applications**

➤ Dynamic Influence Maximization
➤ Approximate Graph Propagation
➤ Computational Epidemiology
➤ Fractional (bipartite) matching

## Technique 1: Group Partition

*Let's start from a simple case!*

$p$ $p$ $p$ $p$ $p$ $p$ $p$
$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ $x_7$

$h = 2$ ...... $i = 3$ → Sample this!

➤ The **index** of the first sample: $i \sim p(1-p)^{i-1}$
➤ The geometric distribution is **memoryless**
➤ The query time = # of the sampled events

*Try a more complicated case!*

➤ $2^{-j} < p(x_i) \leq 2^{-j+1}$
➤ Let $p = 2^{-j+1}$ be the upper bound
➤ First sample each event with $p$ as a **candidate**, then accept it with $\frac{p(x_i)}{p}$
➤ Each event is sampled with probability $p \cdot \frac{p(x_i)}{p} = p(x_i)$
➤ The expect number of candidates $= np \leq 2\mu$ → It costs $O(1 + \mu)$ time

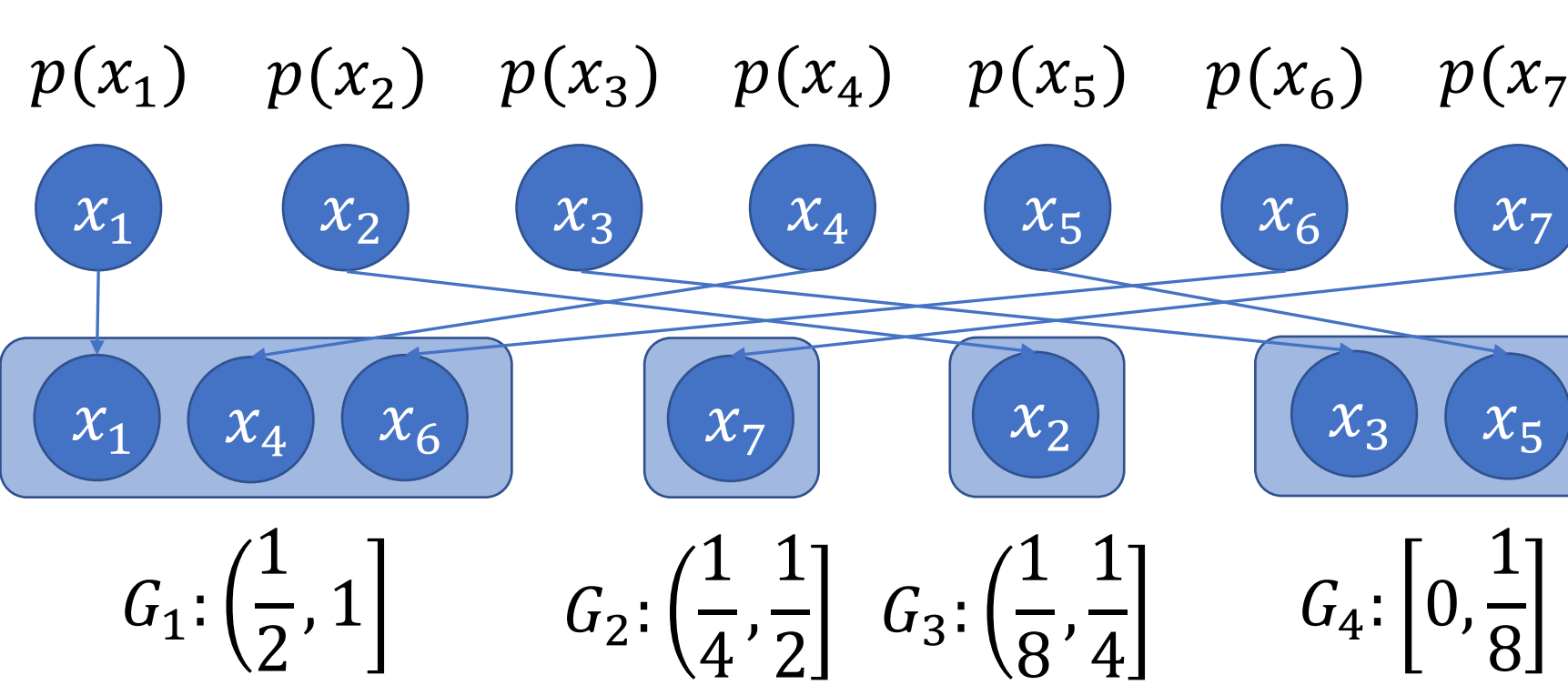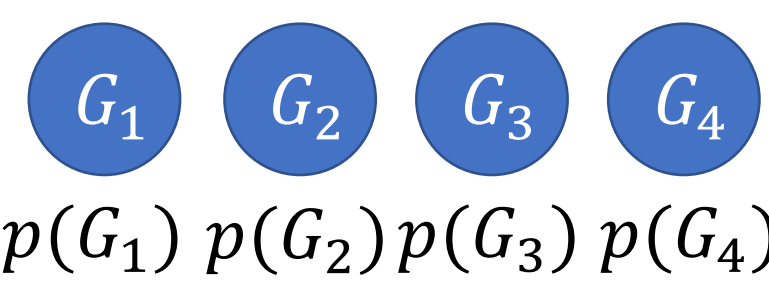### Why Group Partition?

••• *GeoSS*

**Step 0.** Let $p$ be the upper bound
**Step 1.** Currently at the $h$-th event, $h = 0$ initially
**Step 2.** Generate $i \sim p(1-p)^{i-1}$
**Step 3.** The next candidate: $(i + h)$-th event, accept it with $\frac{p(x_i)}{p}$
**Step 4.** $h = i + h$, repeat Step 2 to 4 until $h > n$

**domain set $S$**

$p(x_1)$ $p(x_2)$ $p(x_3)$ $p(x_4)$ $p(x_5)$ $p(x_6)$ $p(x_7)$
$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ $x_7$

$x_1$ $x_4$ $x_6$ | $x_7$ $x_2$ | $x_3$ $x_5$

$G_1: \left(\frac{1}{2}, 1\right]$  $G_2: \left(\frac{1}{4}, \frac{1}{2}\right]$  $G_3: \left(\frac{1}{8}, \frac{1}{4}\right]$  $G_4: \left[0, \frac{1}{8}\right]$

➤ Create ($\lceil \log n \rceil + 1$) groups: $G_1, G_2, \ldots, G_K (K = \lceil \log n \rceil + 1)$
➤ $G_j = \{x_i \mid 2^{-j} < p(x_i) \leq 2^{-j+1}\}$, $1 \leq j \leq K - 1$
➤ $G_j = \{x_i \mid p(x_i) \leq 2^{-j+1}\}, j = K$
➤ Use *GeoSS* within each group
➤ Totally costs $O(1 + \mu + \log n)$ time

### How to $O(1 + \mu + \log n) \to O(1 + \mu)$?

$G_1$ $G_2$ $G_3$ $G_4$

$p(G_1)$ $p(G_2)$ $p(G_3)$ $p(G_4)$

➤ **Only sample the groups with at least one candidate!**
➤ The probability that $G_j$ contains at least one candidate:
$p(G_j) = 1 - \left(1 - 2^{-j+1}\right)^{|G_j|}$
➤ First sample among the groups with $p(G_j)$, then sample within the sampled groups

**Algorithm 1: SampleWithinGroup**
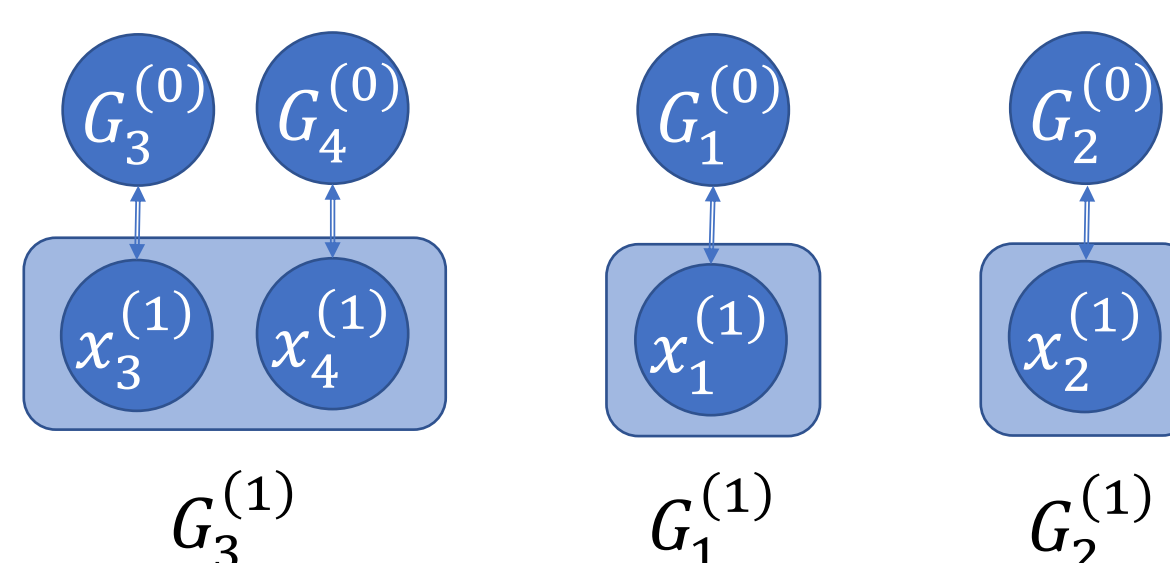
**Input:** a group $G_k$
**Output:** a drawn sample $T$

1 $n_k \leftarrow |G_k|, T \leftarrow \emptyset, h \leftarrow 0$;
2 Let $G_k[i]$ be the $i$-th element of $G_k$;
3 Generate a random $r$ s.t. $\Pr[r = j] = \frac{2^{-k+1}(1-2^{-k+1})^{j-1}}{p(G_k)}$, $j \in \{1, \ldots, n_k\}$;
4 **while** $r + h \leq n_k$ **do**
5   $h \leftarrow r + h$;
6   **if** rand() $< p(G_k[h])/2^{-k+1}$ **then**
7     $T \leftarrow T \cup \{G_k[h]\}$;
8   Generate a random $r \sim$ Geo($2^{-k+1}$);
9 **return** $T$

••• *Tips for updates*

Example: **inserting an event $x$**

**S0.** Insert $x$ to a group $G_k^{(0)}$ based on $p(x)$
**S1.** Recalculate the prob. $p(G_k^{(0)})$
**S2.** Transfer $x_k^{(1)}$ from one group to another according to $p(x_k^{(1)})$ (also $p(G_k^{(0)})$)
**S3.** Recalculate the prob. of the modified groups at **S2**

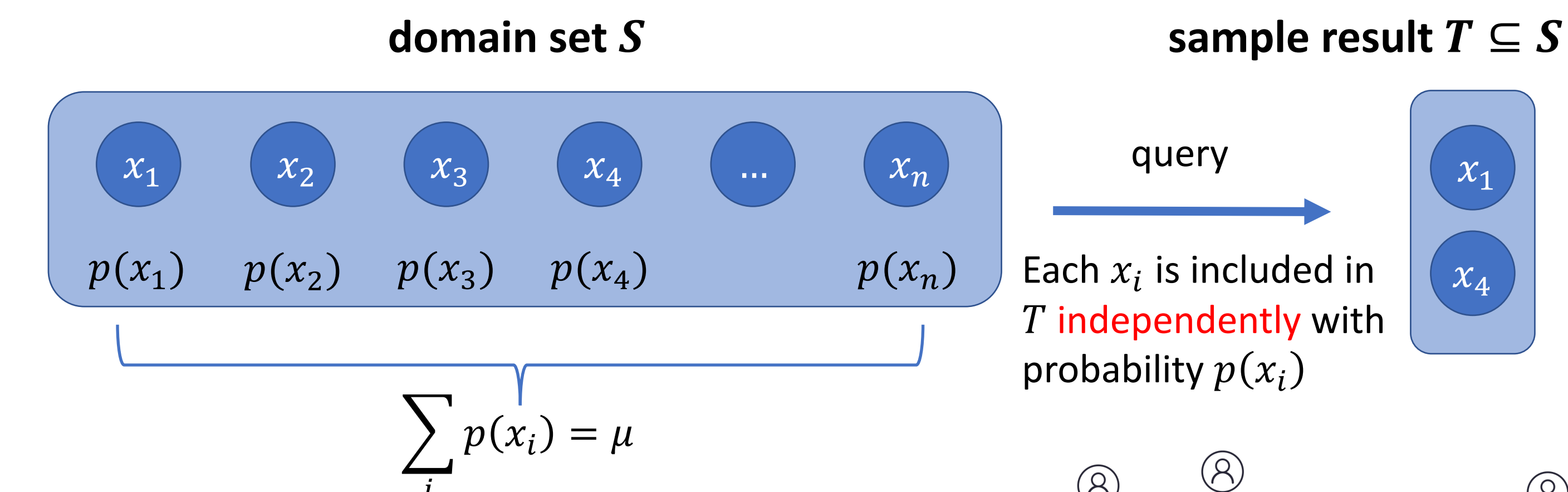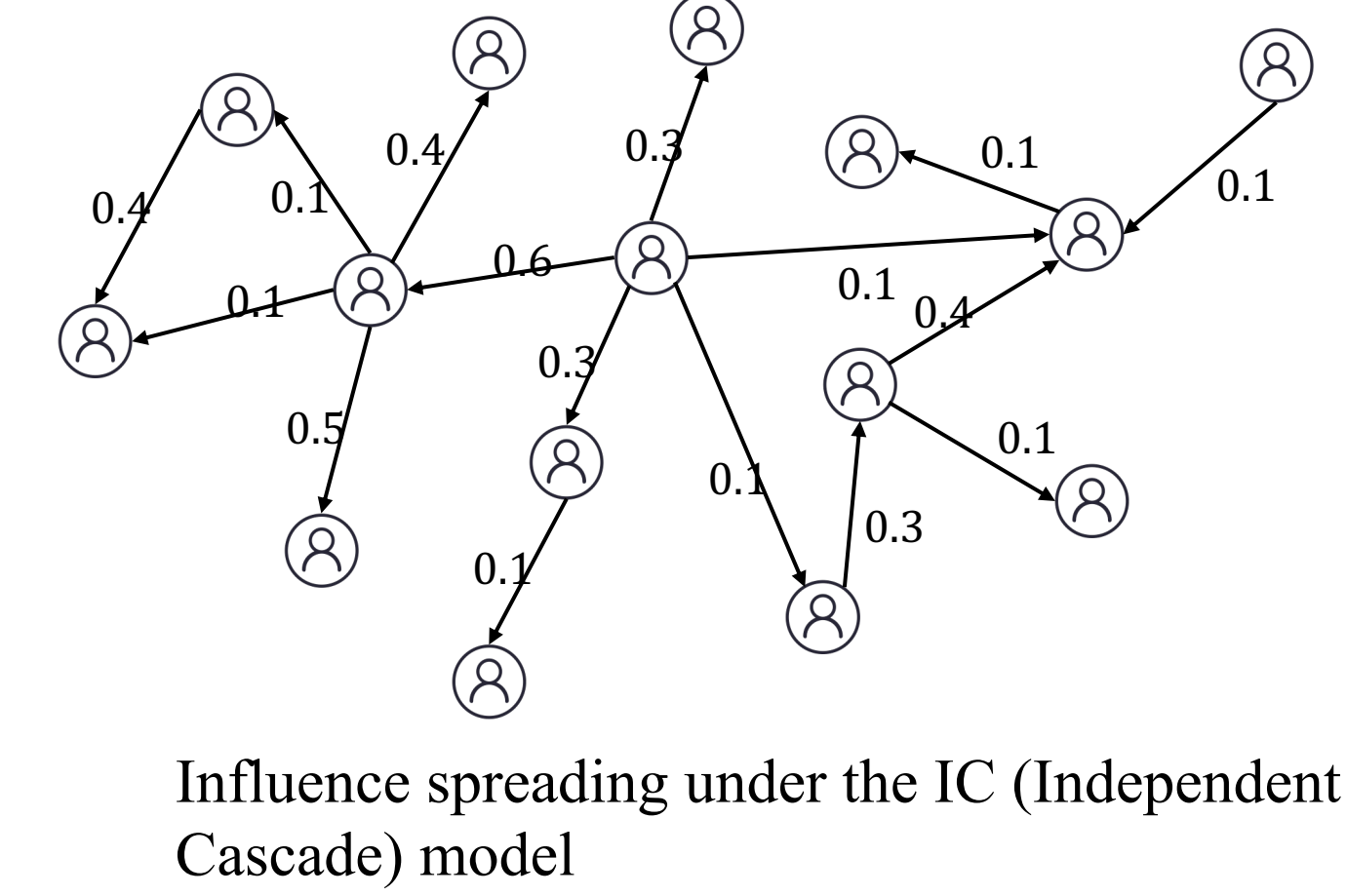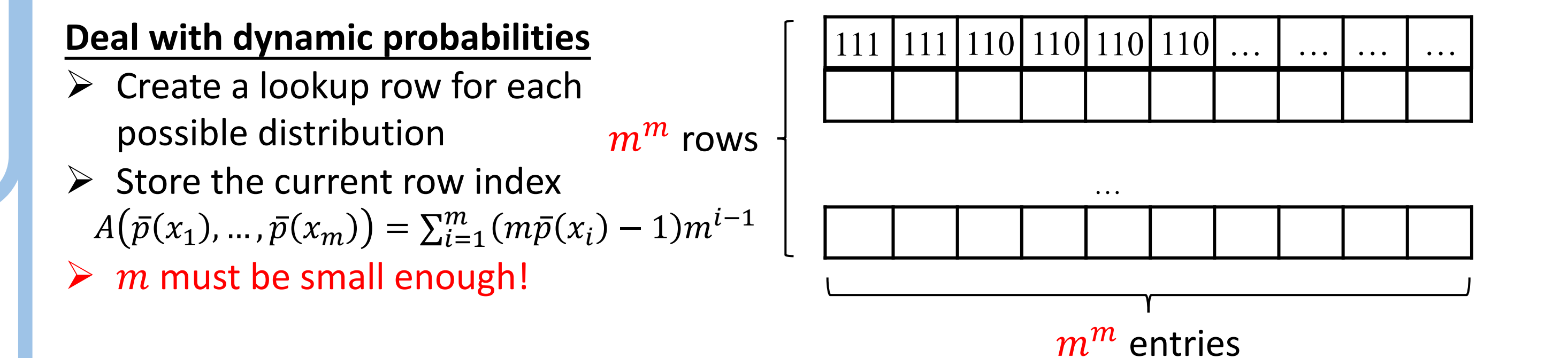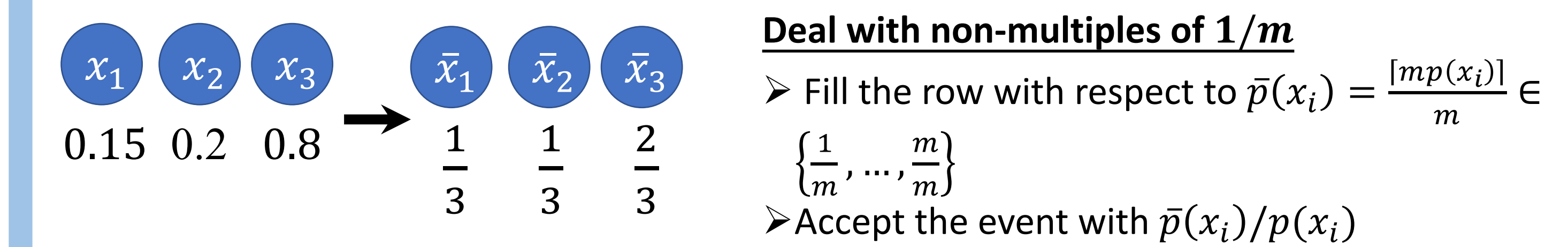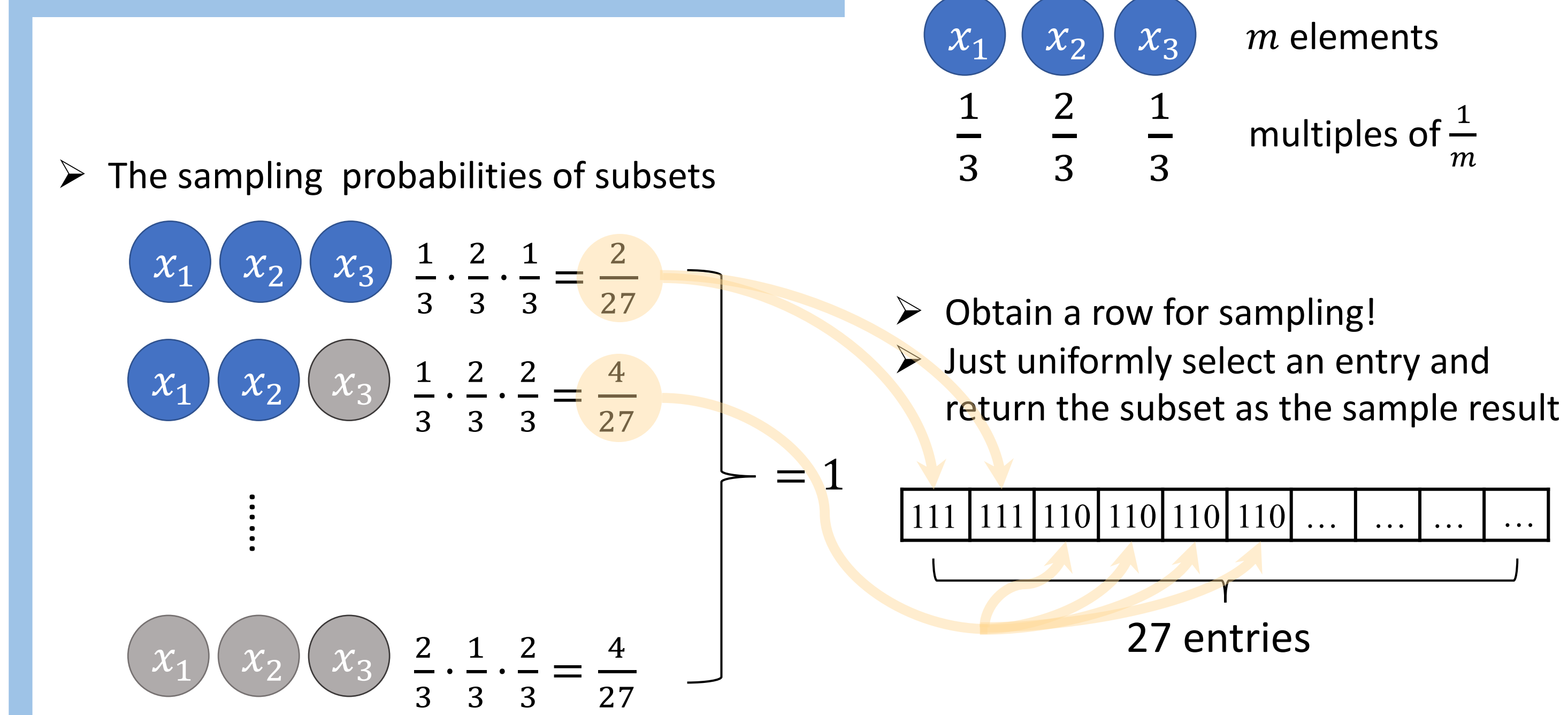### How to sample among the groups?

➤ **Partition again!**
➤ We add level index to distinguish various subset sampling problems
➤ Use **Technique 2** to sample the groups at level 1, only $O(\log \log n)$ events

$G_3^{(0)}$ $G_4^{(0)}$    $G_1^{(0)}$    $G_2^{(0)}$

$x_3^{(1)}$ $x_4^{(1)}$   $x_1^{(1)}$   $x_2^{(1)}$

$G_3^{(1)}$   $G_1^{(1)}$   $G_2^{(1)}$

## Technique 2: Table Lookup

**Sample each element independently ⟺ Sample one subset**

**An example with # of events m=3**

$x_1$ $x_2$ $x_3$   $m$ elements
$\frac{1}{3}$ $\frac{2}{3}$ $\frac{1}{3}$   multiples of $\frac{1}{m}$

➤ The sampling probabilities of subsets

$x_1$ $x_2$ $x_3$   $\frac{1}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} = \frac{2}{27}$
$x_1$ $x_2$ $x_3$   $\frac{1}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} = \frac{4}{27}$
...
$x_1$ $x_2$ $x_3$   $\frac{2}{3} \cdot \frac{1}{3} \cdot \frac{2}{3} = \frac{4}{27}$

$= 1$

➤ Obtain a row for sampling!
➤ Just uniformly select an entry and return the subset as the sample result

111 111 110 110 110 110 ... ... ...

27 entries

$x_1$ $x_2$ $x_3$ → $\bar{x}_1$ $\bar{x}_2$ $\bar{x}_3$
0.15 0.2 0.8    $\frac{1}{3}$ $\frac{1}{3}$ $\frac{2}{3}$

**Deal with non-multiples of $1/m$**

➤ Fill the row with respect to $\bar{p}(x_i) = \frac{\lceil mp(x_i) \rceil}{m} \in \left\{\frac{1}{m}, \ldots, \frac{m}{m}\right\}$
➤ Accept the event with $\bar{p}(x_i)/p(x_i)$

**Deal with dynamic probabilities**

➤ Create a lookup row for each possible distribution
➤ Store the current row index
$A(\bar{p}(x_1), \ldots, \bar{p}(x_m)) = \sum_{i=1}^{m}(m\bar{p}(x_i) - 1)m^{i-1}$
➤ $m$ must be small enough!

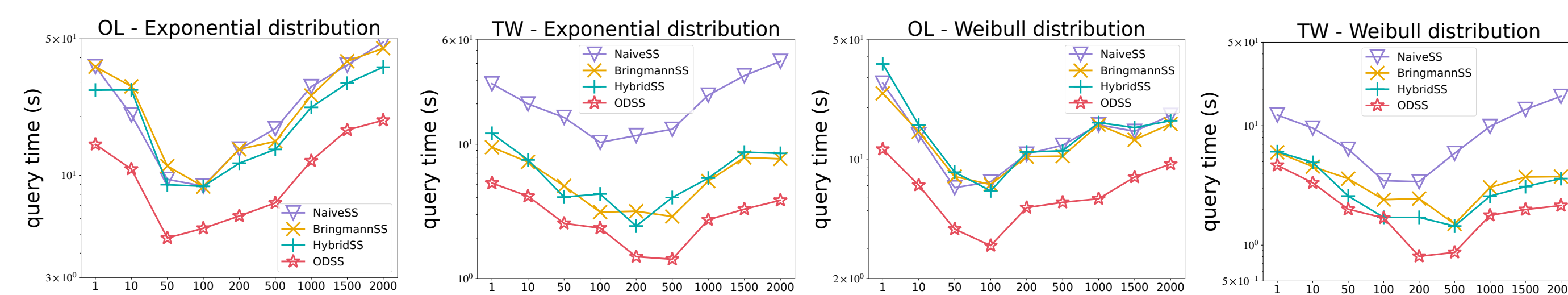111 111 110 110 110 110 ... ...

$m^m$ rows

$m^m$ entries

••• *Tips for updates*

Suppose: $\bar{p}(x_i) \to \bar{p}'(x_i)$, $\bar{p}(x_j) \to \bar{p}'(x_j)$
The new row index:
$A'(\bar{p}(x_1), \ldots, \bar{p}(x_m))$
$= A(\bar{p}(x_1), \ldots, \bar{p}(x_m)) + (m\bar{p}'(x_i) - m\bar{p}(x_i))m^{i-1} + (m\bar{p}'(x_j) - m\bar{p}(x_j))m^{j-1}$

## General Framework

Level 2: $m = \lceil \log(\lceil \log n \rceil + 1) \rceil + 1 = 3$

Level 1: $n^{(1)} = \lceil \log n \rceil + 1 = 4$

Level 0: $n^{(0)} = n = 7$

sample an entry in the $A(\bar{p}(x_1^{(2)}), \bar{p}(x_2^{(2)}), \bar{p}(x_3^{(2)}))$-th row

obtain candidates $x_1^{(2)}, x_2^{(2)}, x_3^{(2)}$ as indicated by 111
accept $x_1^{(2)}, x_2^{(2)}$ and reject $x_3^{(2)}$

row index

0
1
...
111
...
$m^m - 1$

$m^m$ Lookup Table

Group Partition $G_k = \{x_i \mid p(x_i) \in (2^{-k}, 2^{-k+1}]\}$

## Experiments

➤ Competitors

| Algorithm | Expected Query Time | Update Time |
|---|---|---|
| The Naive Method | $O(n)$ | $O(1)$ |
| HybridSS[COCOON'10] | $O\left(1 + n\sqrt{\min\{\bar{p}, 1 - \bar{p}\}}\right)$ | $O(n)$ |
| BringmannSS[ICALP'12] | $O(1 + \mu)$ | $O(\log^2 n)$ |
| ODSS (Ours) | $O(1 + \mu)$ | $O(1)$ |

➤ Distributions of probabilities
• Normal distribution (skewness as 0)
• Half-normal distribution (skewness below 1)
• Exponential distribution (skewness as 2)
• Log-normal distribution (skewness as 4)
• Re-scale the range of the random number into [0,1]

➤ Trade-off between query time and update time. ($n = 10^5, \mu = 1$)

➤ Varying $\mu$: query time (s) on distributions with different skewnesses. ($n = 10^6$)

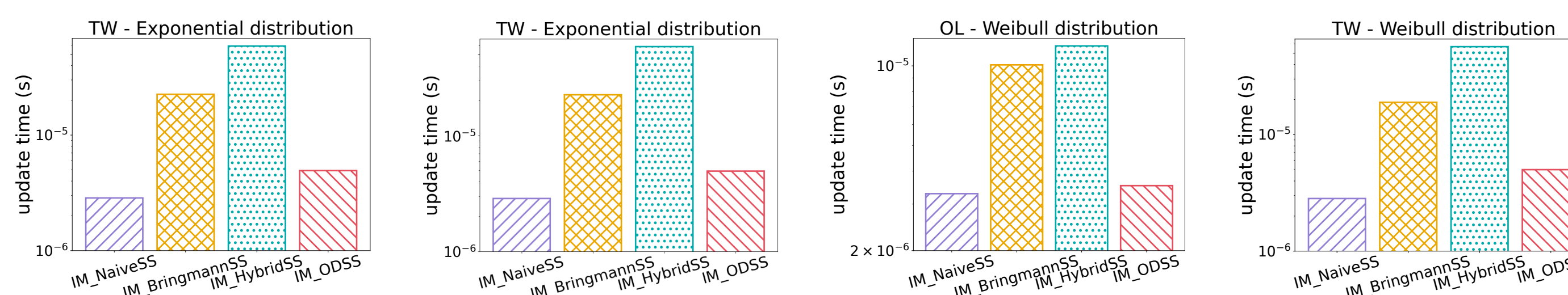➤ Varying $n$: update time (s) on distributions with different skewnesses.

## Empirical Study on Influence Maximization

➤ Based on the framework OPIM-C[ICMD'18], replace the subset sampling module with various dynamic subset sampling structures and thus obtain a new dynamic IM algorithm for the fully dynamic model.
➤ No algorithms can achieve any meaningful approximation guarantee in the fully dynamic network model. That is, re-running an IM algorithm upon each update can achieve the lower bound of the running time.

OL - Exponential distribution | TW - Exponential distribution | OL - Weibull distribution | TW - Weibull distribution

➤ Running time of dynamic IM algorithms based on various subset sampling structures.

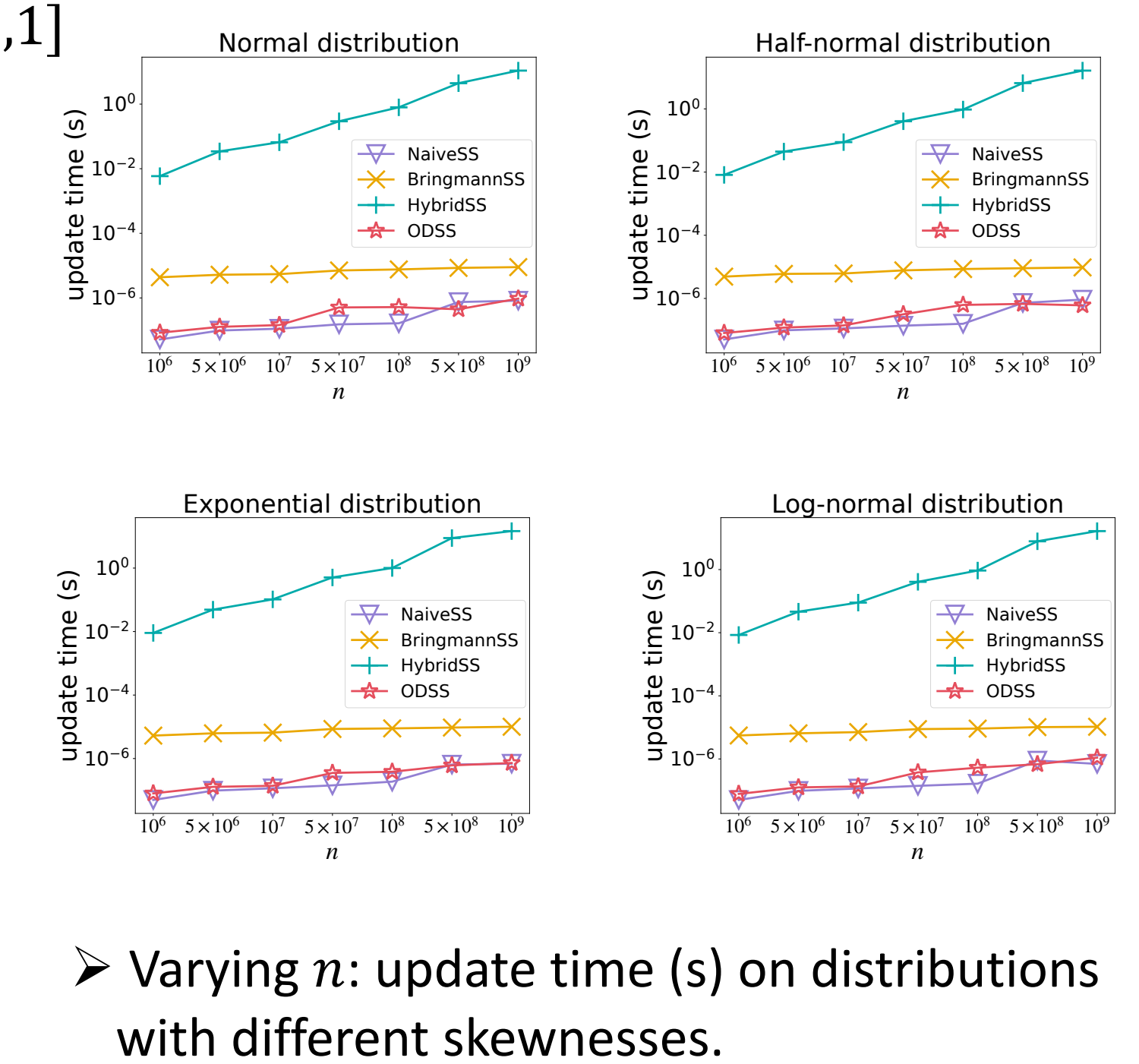TW - Exponential distribution | TW - Exponential distribution | OL - Weibull distribution | TW - Weibull distribution

➤ Update time of dynamic IM algorithms based on various subset sampling structures.