

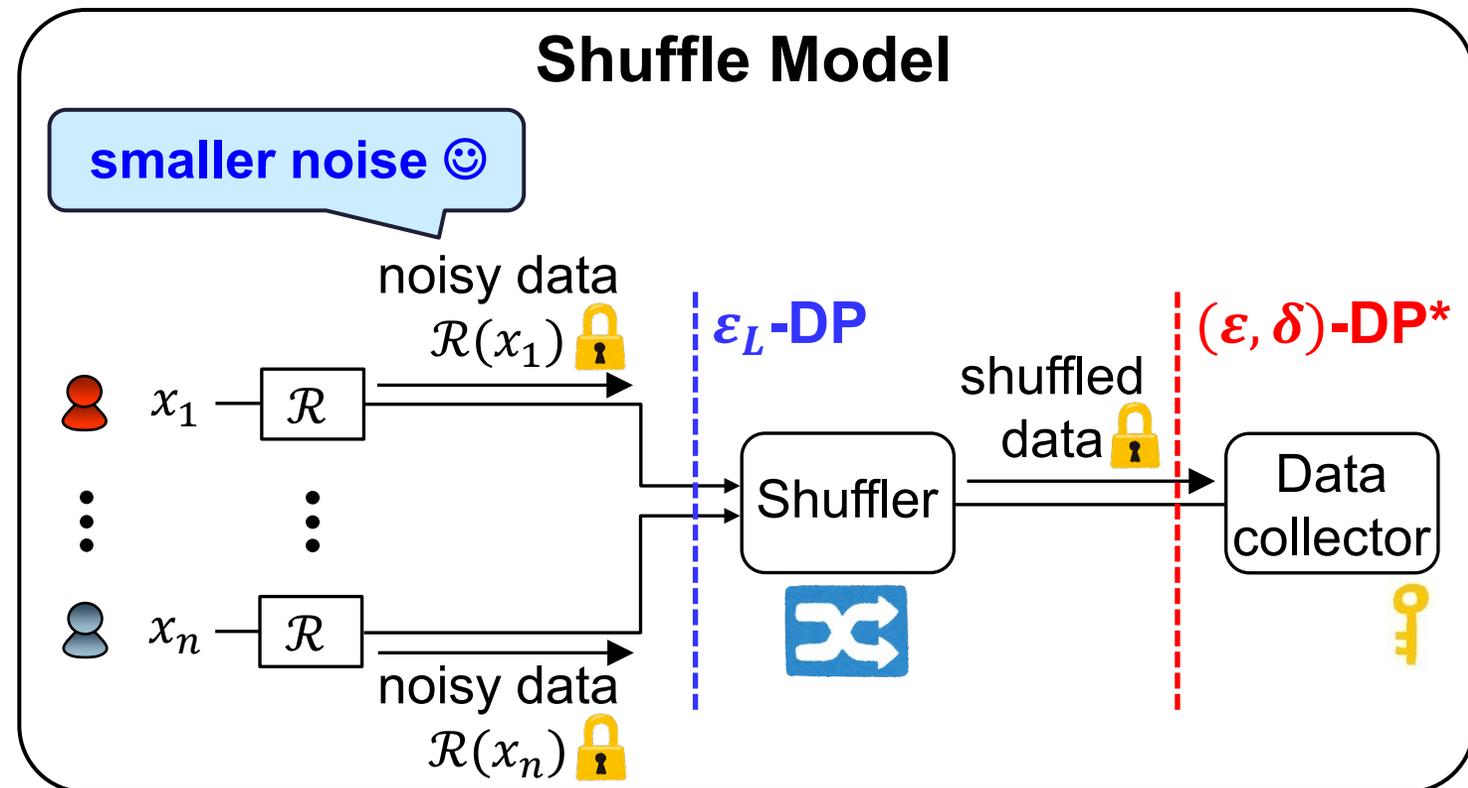
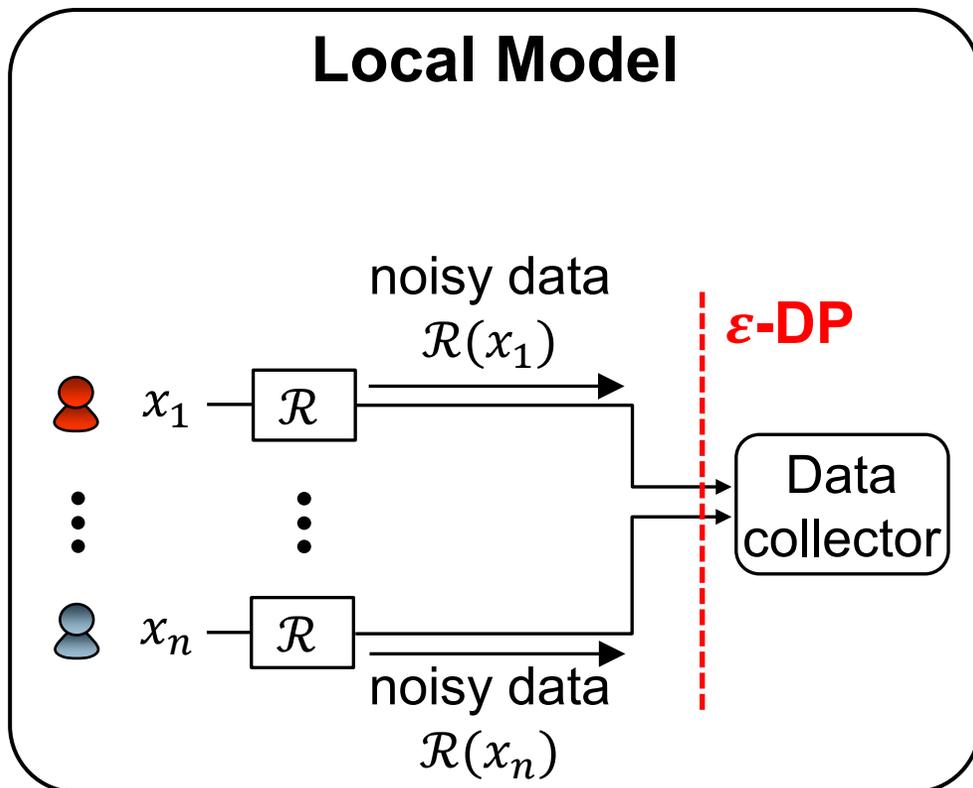
Augmented Shuffle Differential Privacy Protocols for Large-Domain Categorical and Key-Value Data

Takao Murakami (ISM/AIST/RIKEN AIP), Yuichi Sei (UEC), Reo Eriguchi (AIST)

Shuffle DP (Differential Privacy)

▶ Shuffle Model

- ▶ Introduces a *shuffler*, which does not collude with a data collector.
- ▶ Shuffling amplifies privacy ($\epsilon \ll \epsilon_L$). \rightarrow Accuracy is improved at the same value of ϵ .



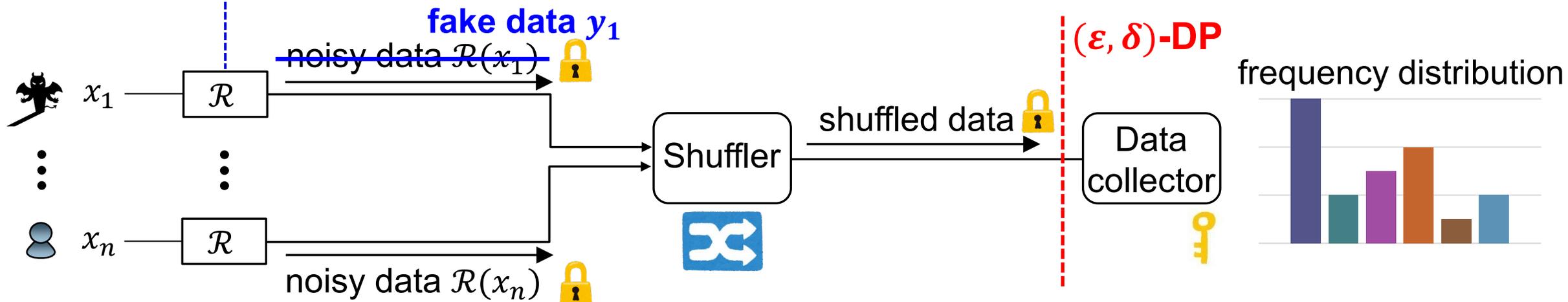
* δ is extremely small (e.g., $\delta \ll 1/n$).

Two Issues in Shuffle DP

▶ Issue 1: Data Poisoning Attacks

- ▶ Some malicious users may send fake data to manipulate the statistics about input data.
- ▶ When ϵ is close to 0, normal users (still) have to add large noise, but malicious users do not.
- ▶ → Accuracy is significantly degraded by poisoning attacks.

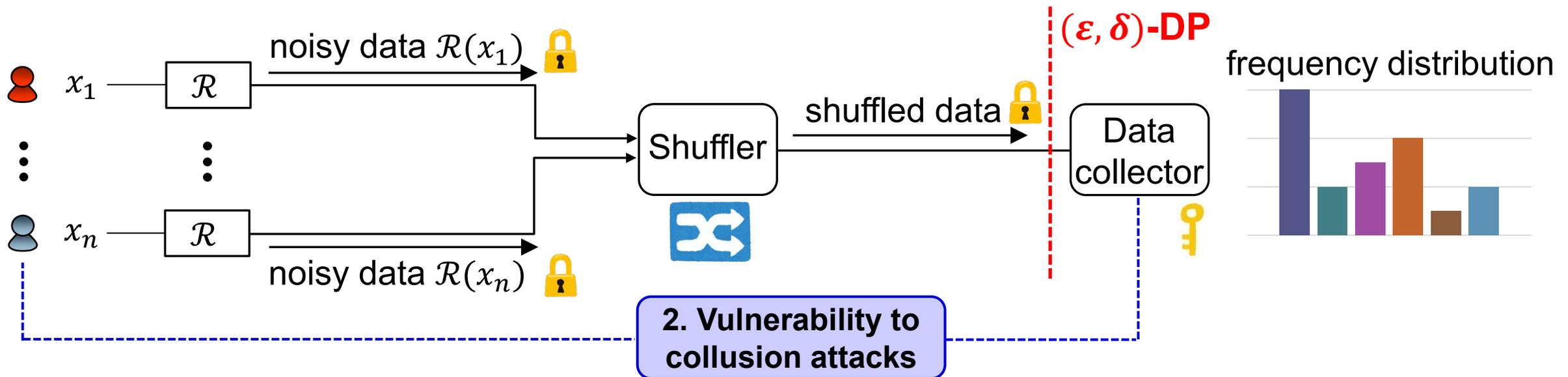
1. Vulnerability to data poisoning attacks



Two Issues in Shuffle DP

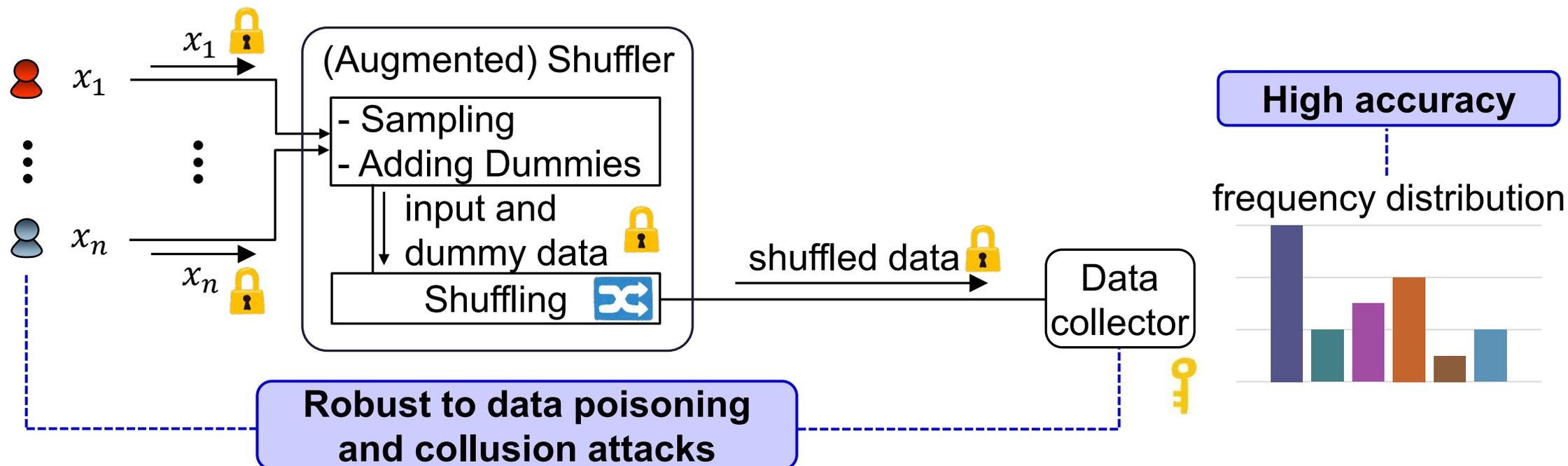
▶ Issue 2: Collusion Attacks

- ▶ Some users may share their noisy data with the data collector to reduce the shuffling effect.
 - ▶ E.g. When $n - 1$ users collude with the data collector, we will get no shuffling effect.
 - ▶ E.g. When #colluding users is $0.1n$ ($n = 6 \times 10^5$, $\delta = 10^{-12}$), ϵ is increased from 1 to 7.2.
- ▶ This is an issue, as attackers can inject many malicious accounts in practice [Thomas+, SEC13].



Augmented Shuffle DP

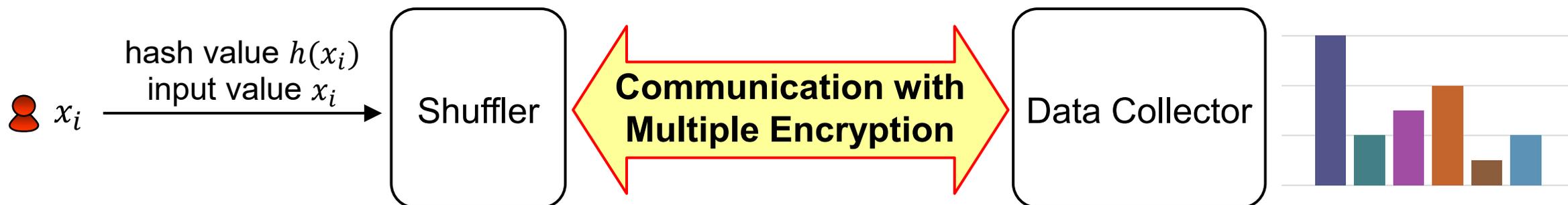
- ▶ LNF (Local-Noise-Free) Protocol [Murakami+, S&P25]
 - ▶ Introduces additional operations (e.g., sampling, adding dummies) into the shuffler.
 - ▶ Achieves high accuracy and robustness against data poisoning and collusion attacks. 😊



- ▶ Cannot be applied to large-domain data due to high communication/computational costs, e.g., it requires 3 years when $d = 10^9$ (d : #items). 😞

This Work

- ▶ Our Proposal: FME (Filtering-with-Multiple-Encryption) Protocol
 - ▶ Improves the efficiency of LNF by carefully using hashing and **multiple encryption**. 😊
 - ▶ To our knowledge, we are the first to use multiple encryption in the DP literature.
 - ▶ We applied FME to frequency estimation and KV (Key-Value) statistics estimation.



- ▶ We show that FME provides (computational) DP and significantly improves the efficiency:
 - ▶ 3 years \rightarrow 1 day when $d = 10^9$ (d : #items).
 - ▶ d is smaller than 10^9 in most practical applications, e.g., $d = 6 \times 10^8$ in Amazon.

Contents

LNf (Local-Noise-Free) Protocol
[Murakami+, S&P25]

FME (Filtering-with-Multiple-Encryption)
Protocol [This Work]

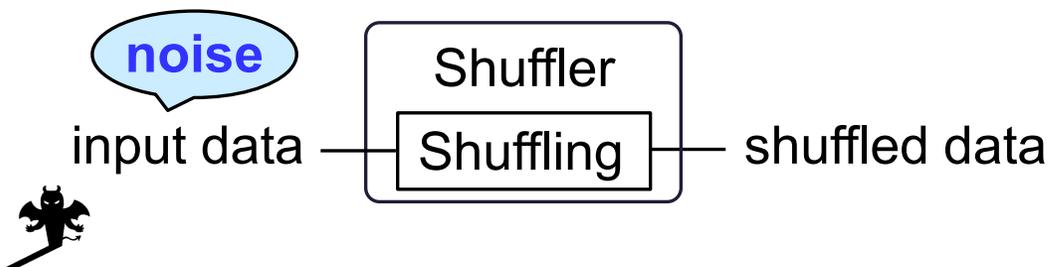
Conclusions

LNf Protocol

► Overview

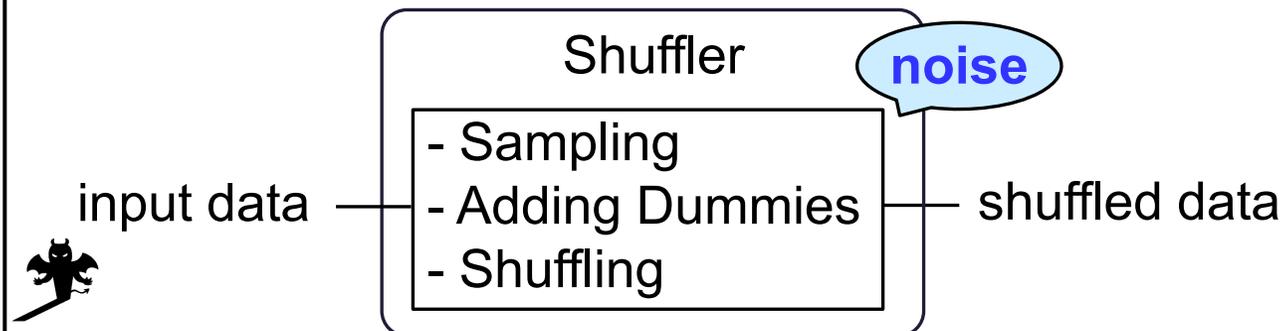
- In most existing protocols, users add noise to their data. → Vulnerable to poisoning and collusion.
- LNf prevents the malicious users' attacks by adding noise **on the shuffler side**.
- Shuffler is (still) simple and can be implemented with any PKE scheme (e.g. RSA, ECIES).

Most Existing Protocols (Pure Shuffle Model)



Vulnerable to poisoning and collusion ☹️

LNf Protocol (Augmented Shuffle Model)

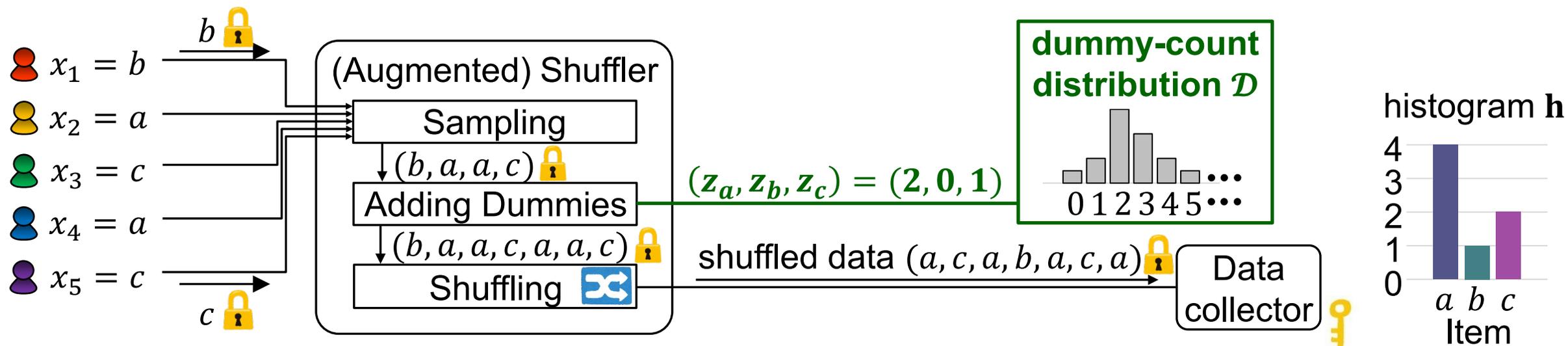


Robust against poisoning and collusion 😊

LNf Protocol

LNf Protocol $\mathcal{S}_{\mathcal{D},\beta}$

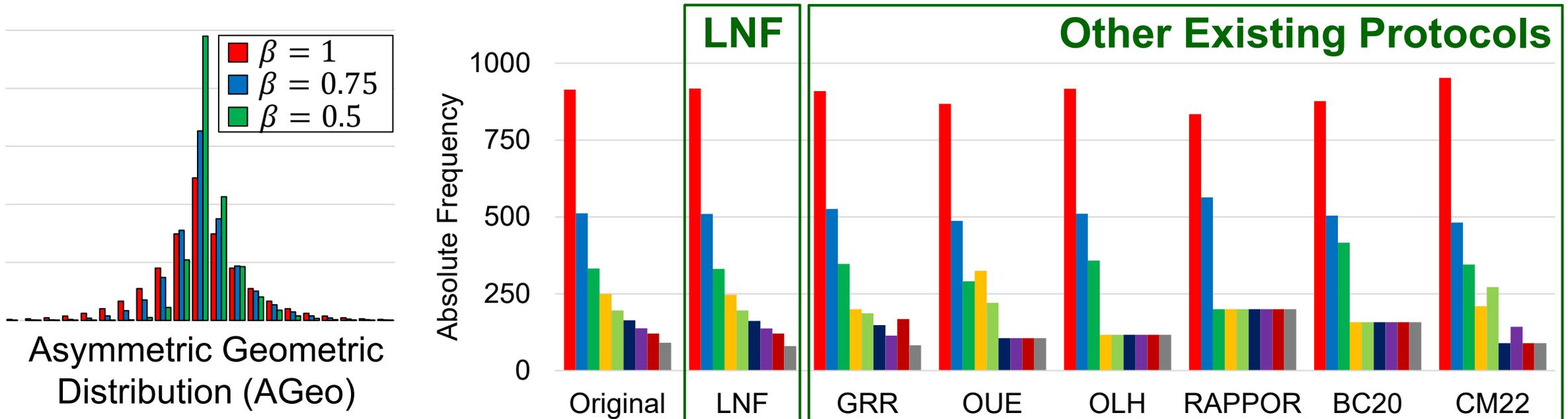
1. Users send their (encrypted) input data without adding noise.
2. Shuffler performs the following three operations:
 - ▶ **Sampling:** Sample each input value with probability $\beta \in [0,1]$.
 - ▶ **Adding Dummies:** Add $z_i \sim \mathcal{D}$ (encrypted) dummies for each item i .
 - ▶ **Shuffling:** Shuffle input and dummy data.
3. Data collector calculates an unbiased estimate of the frequency distribution from \mathbf{h} .



$\mathcal{S}_{\mathcal{D},\beta}$ provides DP and is robust against poisoning and collusion attacks if \mathcal{D} provides DP.

Accuracy and Efficiency

- ▶ Accuracy (Census Dataset, $\varepsilon = 1$)
 - ▶ LNF uses AGeo as the dummy-count distribution \mathcal{D} providing DP.
 - ▶ LNF is very accurate even for unpopular items. 😊



- ▶ Efficiency
 - ▶ LNF suffers from high communication/computational costs $O(d)$ (d : #items) and cannot be applied to large-domain data; e.g., 100 Tbits and 3 years when $d = 10^9$. 😞

Contents

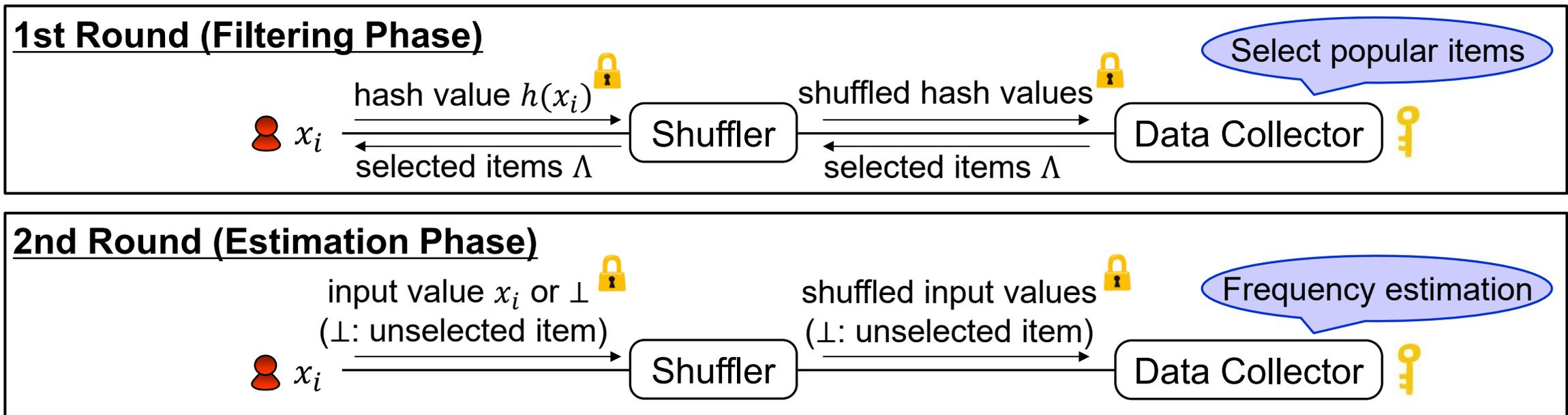
LNf (Local-Noise-Free) Protocol
[Murakami+, S&P25]

**FME (Filtering-with-Multiple-Encryption)
Protocol [This Work]**

Conclusions

Technical Motivation

- ▶ How to Reduce the Domain Size?
 - ▶ We can reduce the domain size by using a hash function $h: [d] \rightarrow [b]$ ($b \ll d$).
 - ▶ Hash values cannot be directly used for frequency estimation due to hash collision, but can be used for filtering items, i.e., selecting popular items.
- ▶ Two-Round Protocol (Strawman Approach)

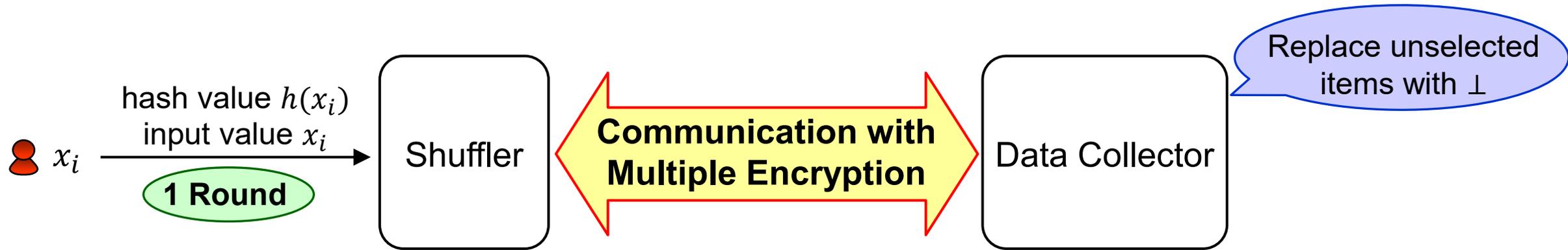


[Problem]: “2-round” interaction greatly reduces usability. ☹️

FME (Filtering-with-Multiple-Encryption) Protocol

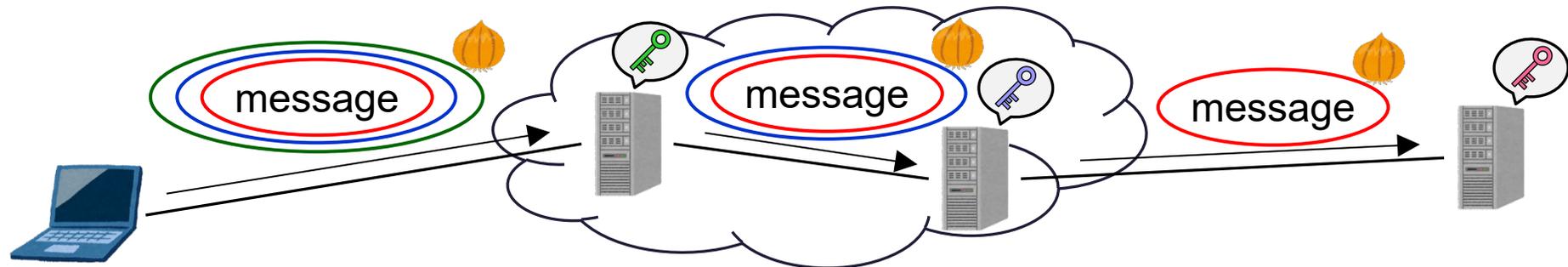
▶ Our Key Idea

- ▶ Reduce #rounds by replacing unselected items with \perp **on the data collector side**.
- ▶ To do this while preserving DP, we carefully use **multiple encryption** between two servers.



▶ Multiple Encryption

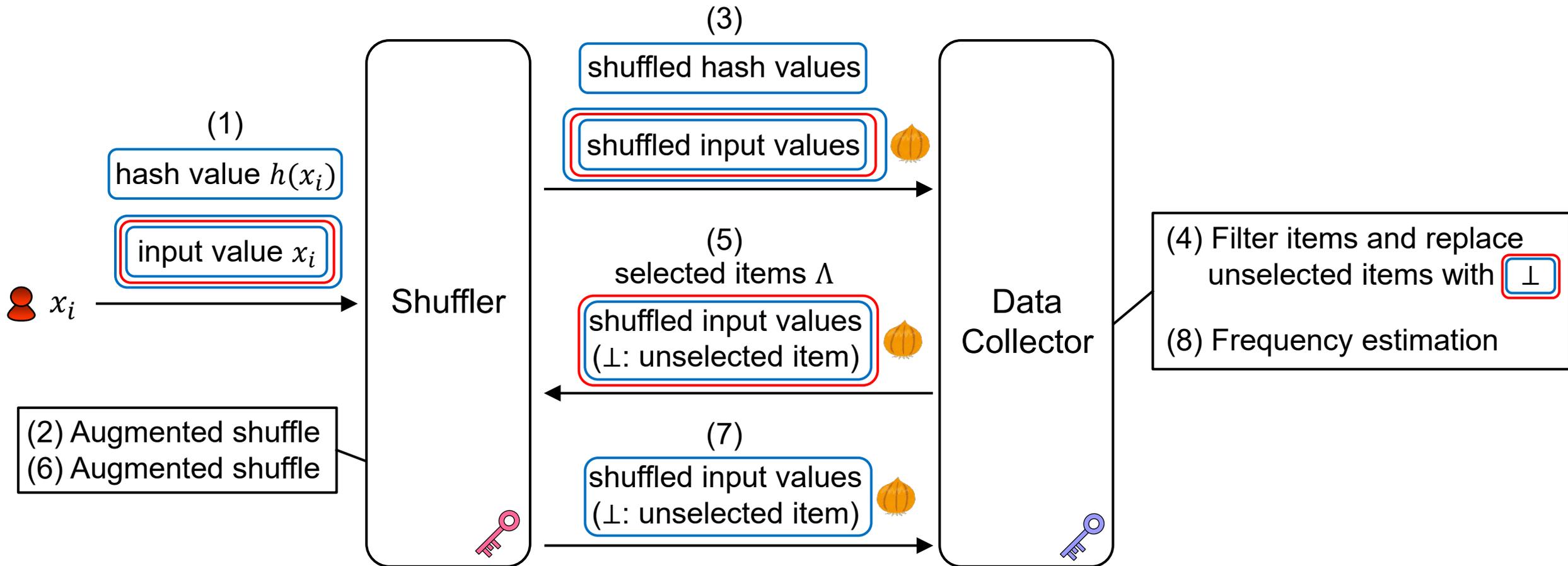
- ▶ Used for anonymous communication (onion routing).
- ▶ We are the first to use this technique (+ “ \perp replacement” trick) to reduce #rounds under DP.



FME (Filtering-with-Multiple-Encryption) Protocol

Protocol

- We use multiple encryption to make “shuffled input values” in steps (3)(5)(7) completely different from each other. → Prevent any attack against these data.



We proved that our FME protocol provides computational DP (→ our paper).

Applying FME to KV (Key-Value) Data

▶ KV Data

- ▶ Each user has key-value pairs (e.g., $\langle \text{Star Wars}, 4 \rangle, \langle \text{Godfather}, 5 \rangle$)
- ▶ Data collector estimates the frequency and mean for each item.

▶ Our Approach (Overview)

1. Discretize the value to ± 1 using padding-and-sampling [Gu+, SEC20].
2. Transform KV pairs into one-dim data ($2d$ categories) and filter the data *at a key level*.

Filtering at a KV pair level (strawman)

		Key					
		1	2	3	4	...	d
Value	1	4	5	0	0	...	1
	-1	3	2	1	0	...	4

: selected



Large bias for the mean value ☹️

Filtering at a key level (our approach)

		Key					
		1	2	3	4	...	d
Value	1	4	5	0	0	...	1
	-1	3	2	1	0	...	4

: selected

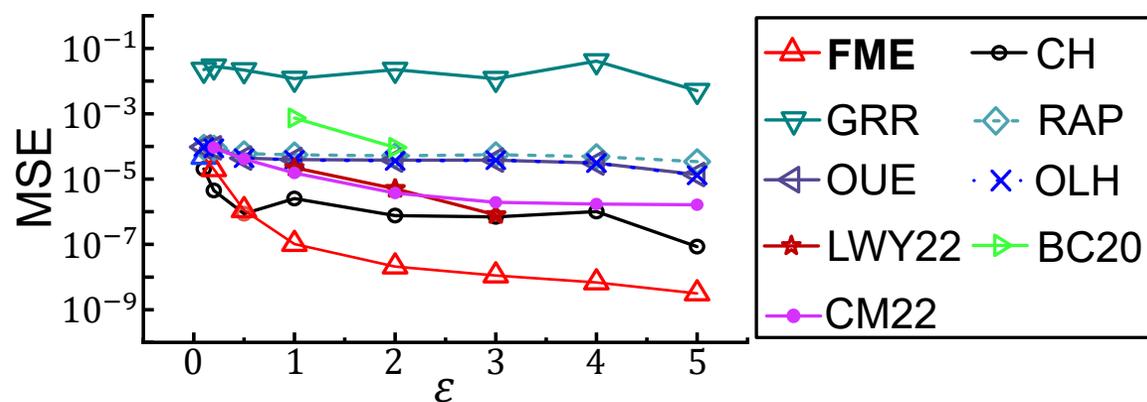


Reduce bias 😊

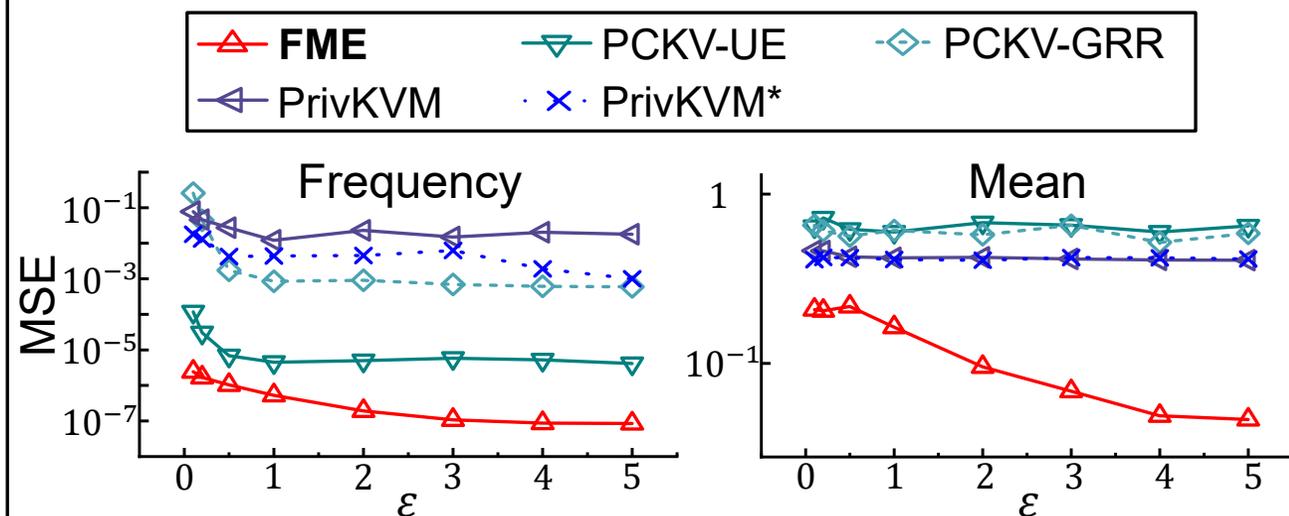
Experimental Results

- ▶ MSE ($\delta = 10^{-12}$)
 - ▶ Our FME significantly outperforms the existing protocols.

Frequency Estimation (Foursquare, $d = 10^6$)



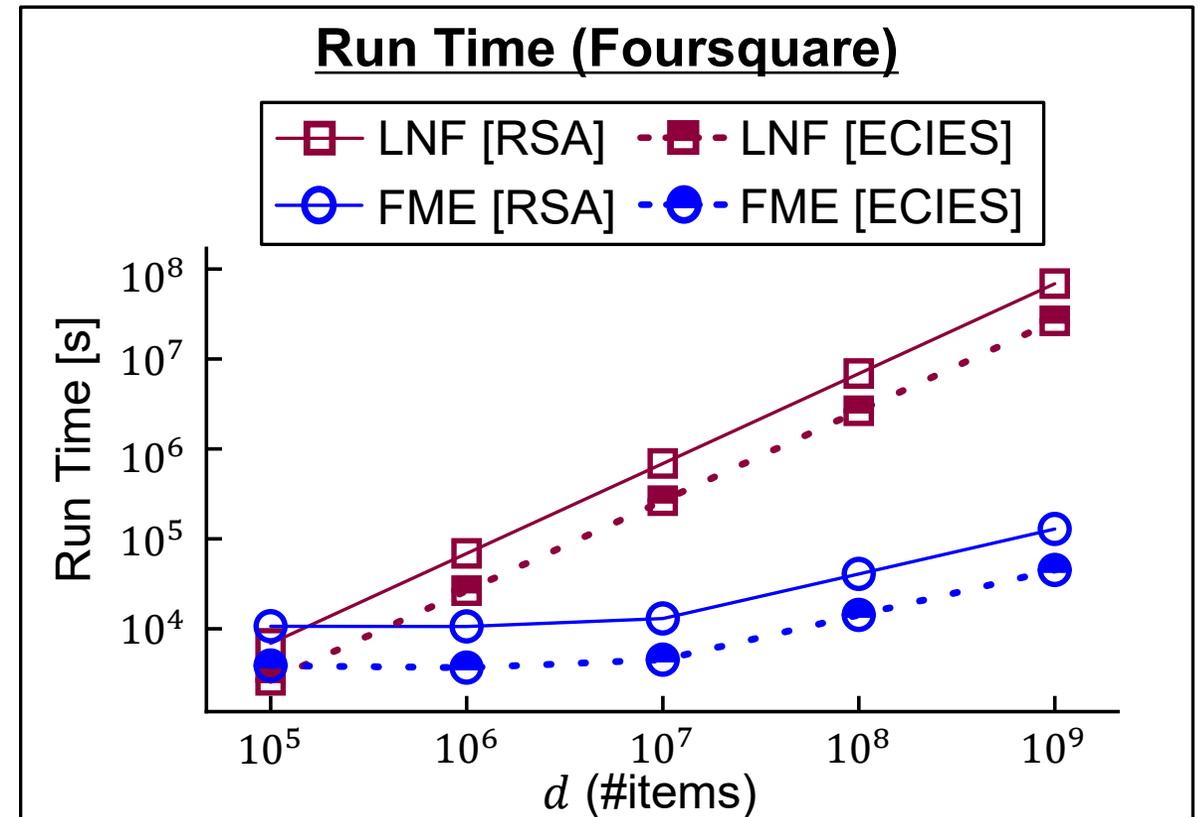
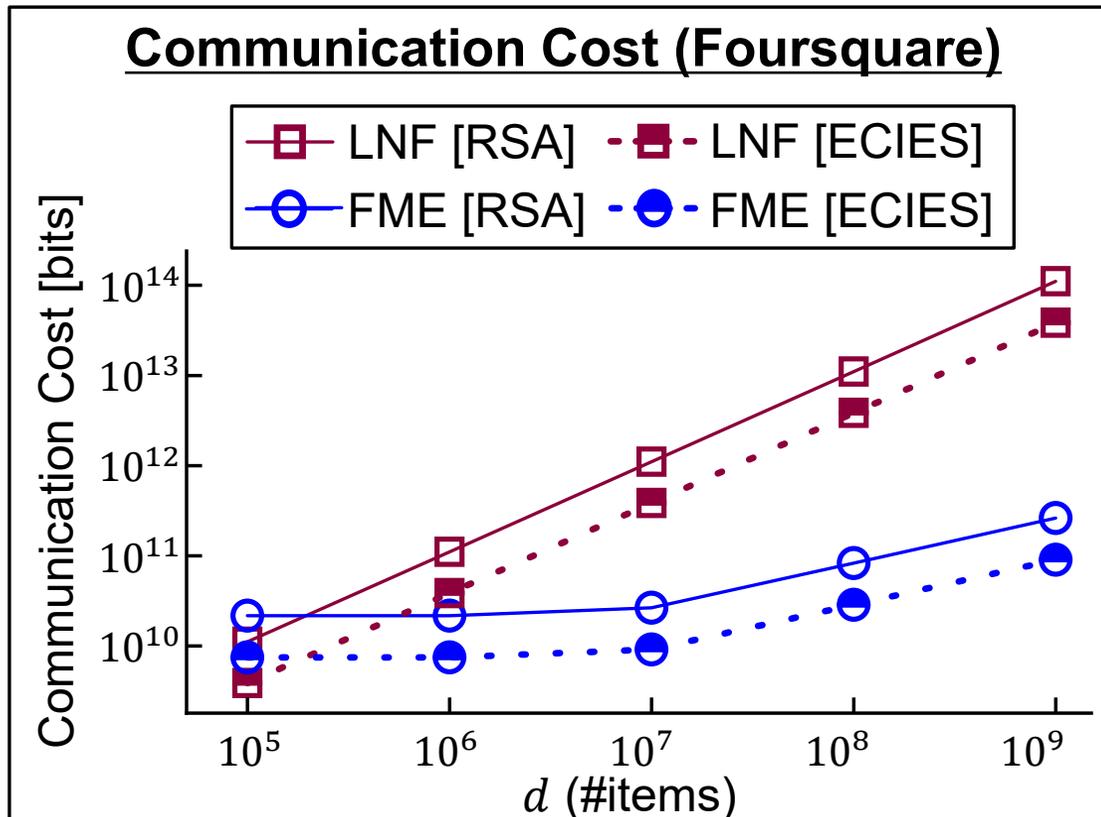
KV Statistics Estimation (Amazon, $d = 2 \times 10^5$)



Experimental Results

► Efficiency

- Our FME can reduce the cost from $O(d)$ to $O(\sqrt{d})$ w/o significantly affecting accuracy.
- When $d = 10^9$, FME reduces the cost from 100 Tbits to 260 Gbits and 3 years to 1 day.



We also showed the robustness of FME against poisoning and collusion attacks.

Contents

LNf (Local-Noise-Free) Protocol
[Murakami+, S&P25]

FME (Filtering-with-Multiple-Encryption)
Protocol [This Work]

Conclusions

Conclusions

▶ Summary

- ▶ We proposed the FME protocol, which significantly improves the efficiency of LNF (e.g., 3 years → 1 day) by using hashing and multiple encryption.

▶ Future Work

- ▶ Applying our protocols to other tasks, e.g., frequent itemset mining, ranking estimation.

Thank you for your attention!

Contact: tmura@ism.ac.jp