

AINA2023

Expectation-Maximization Estimation for Key-Value Data Randomized with Local Differential Privacy

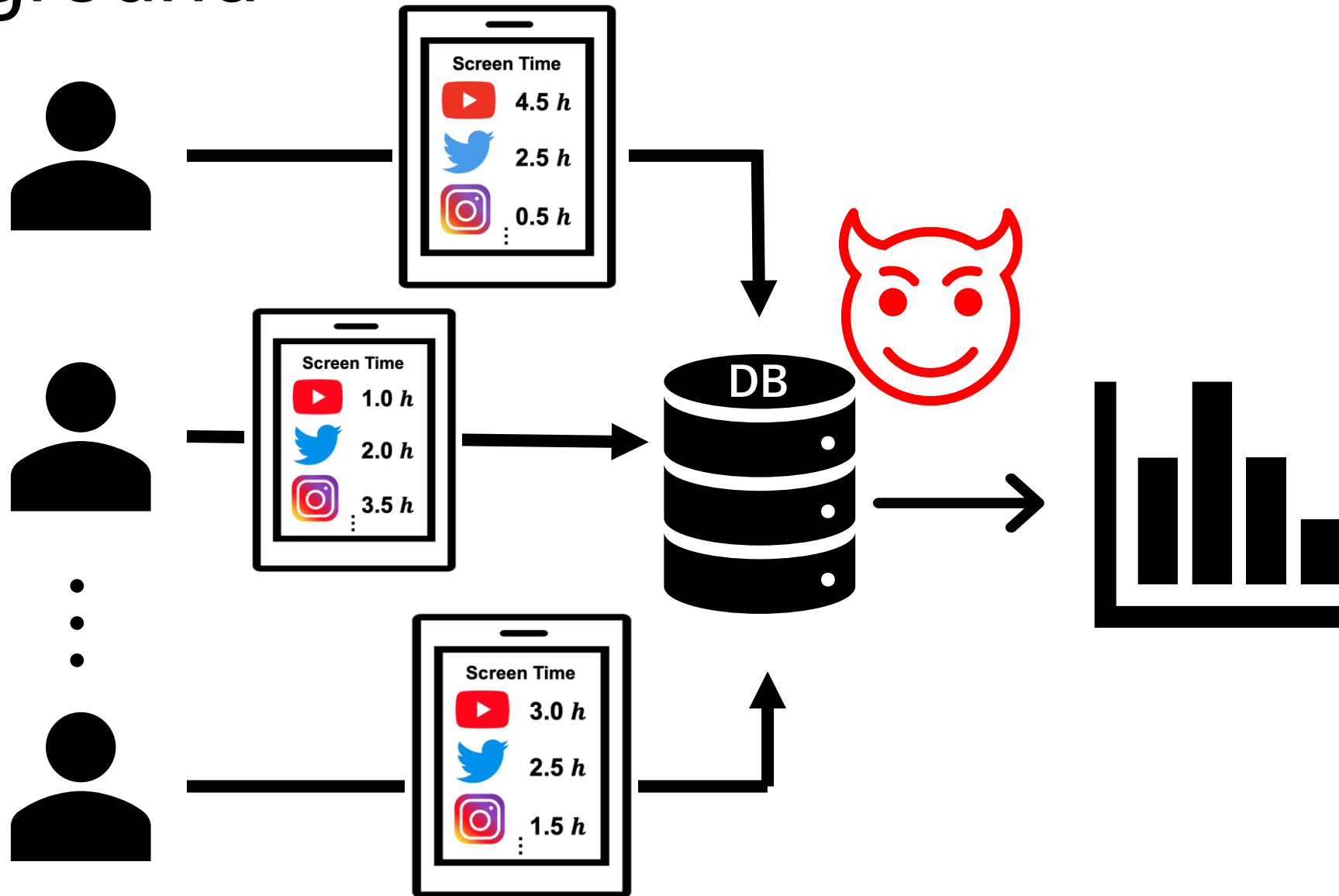
Hikaru Horigome, Hiroaki Kikuchi

(Meiji University)

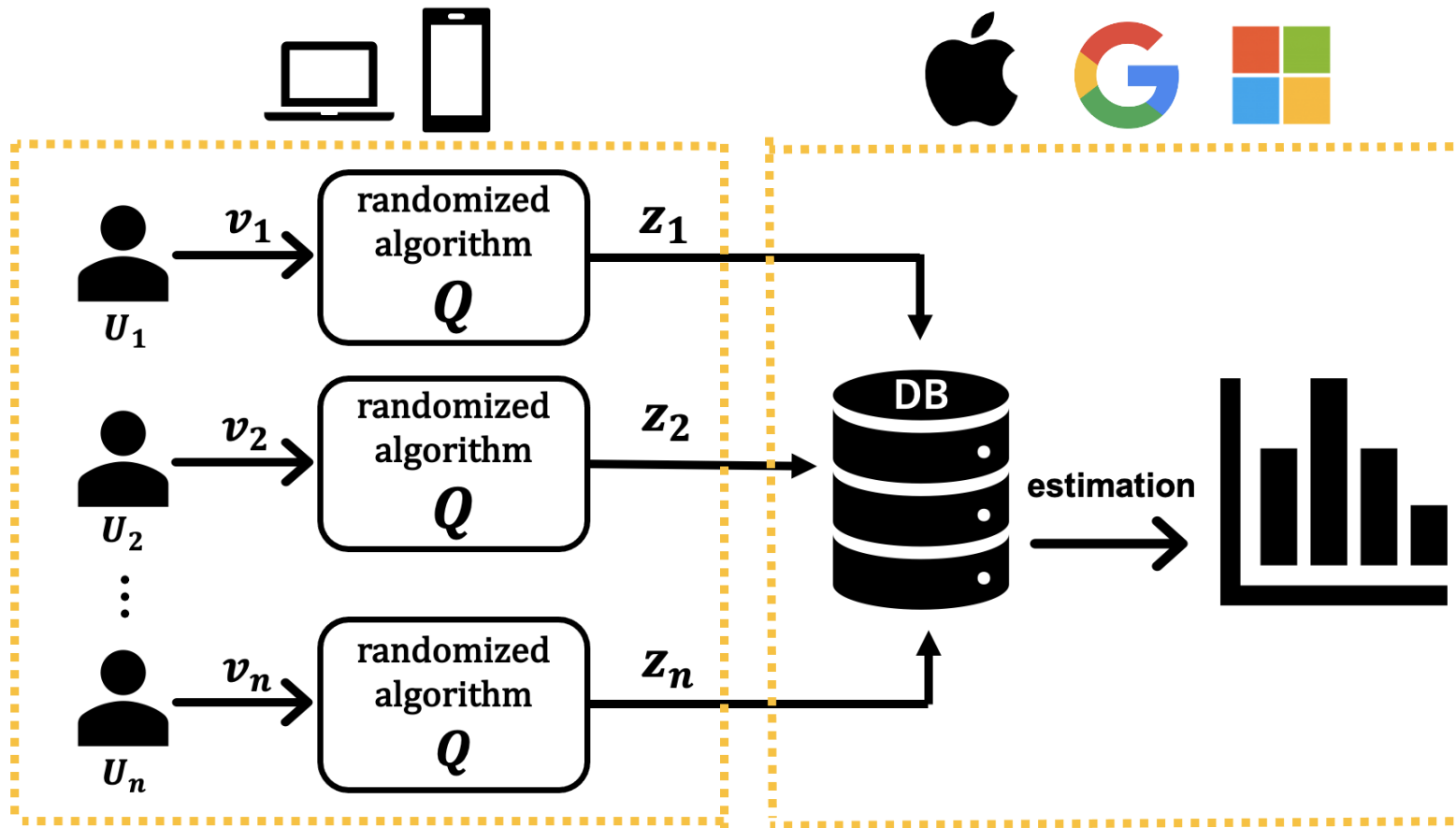
Chia-Mu Yu

(National Yang Ming Chiao Tung University)

Background



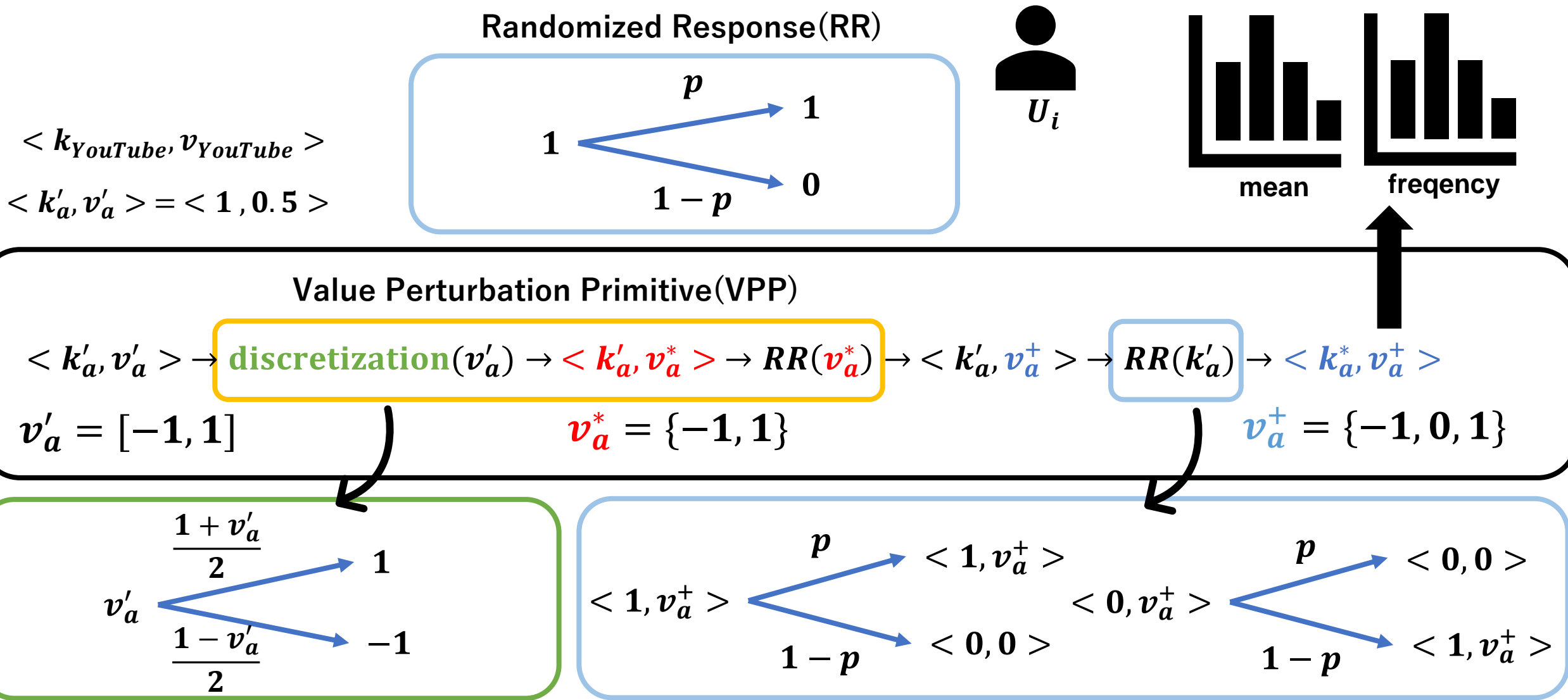
Local Differential Privacy (LDP) [Duchi, et al., 2013]



A randomized algorithm Q satisfies ϵ -LDP if for all pair of values $v, v' \in V$ and all subset S of range Z ($S \subset Z$)

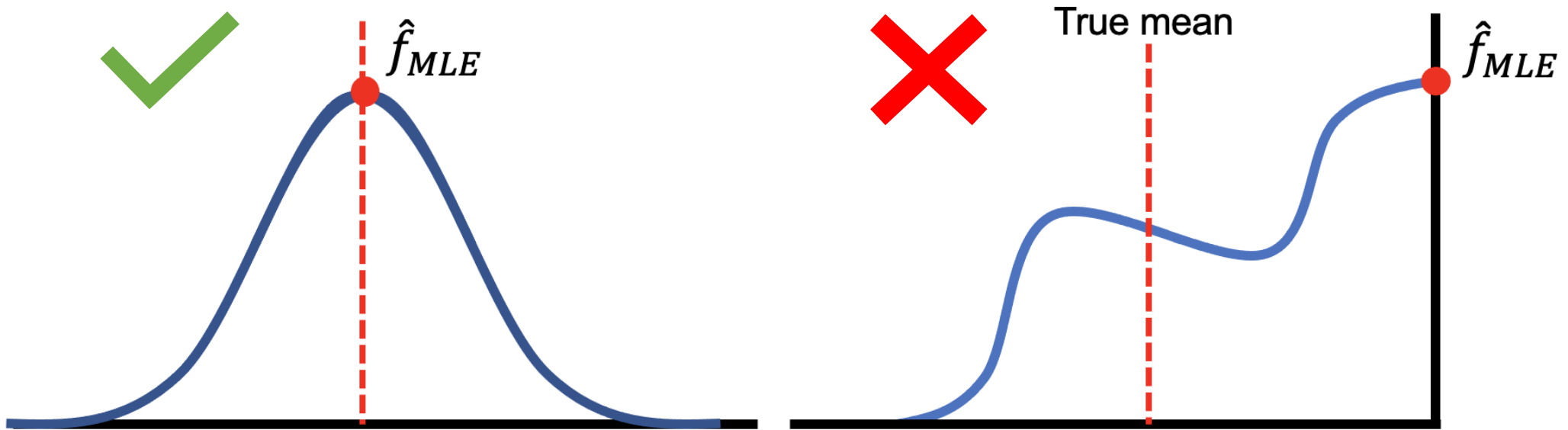
$$\frac{\Pr[Q(v) \in S]}{\Pr[Q(v') \in S]} \leq e^\epsilon$$

PrivKV [Q.Ye , et al., 2019]



Problems of PrivKV

- Most Likelihood Estimation(MLE) : Low estimation accuracy



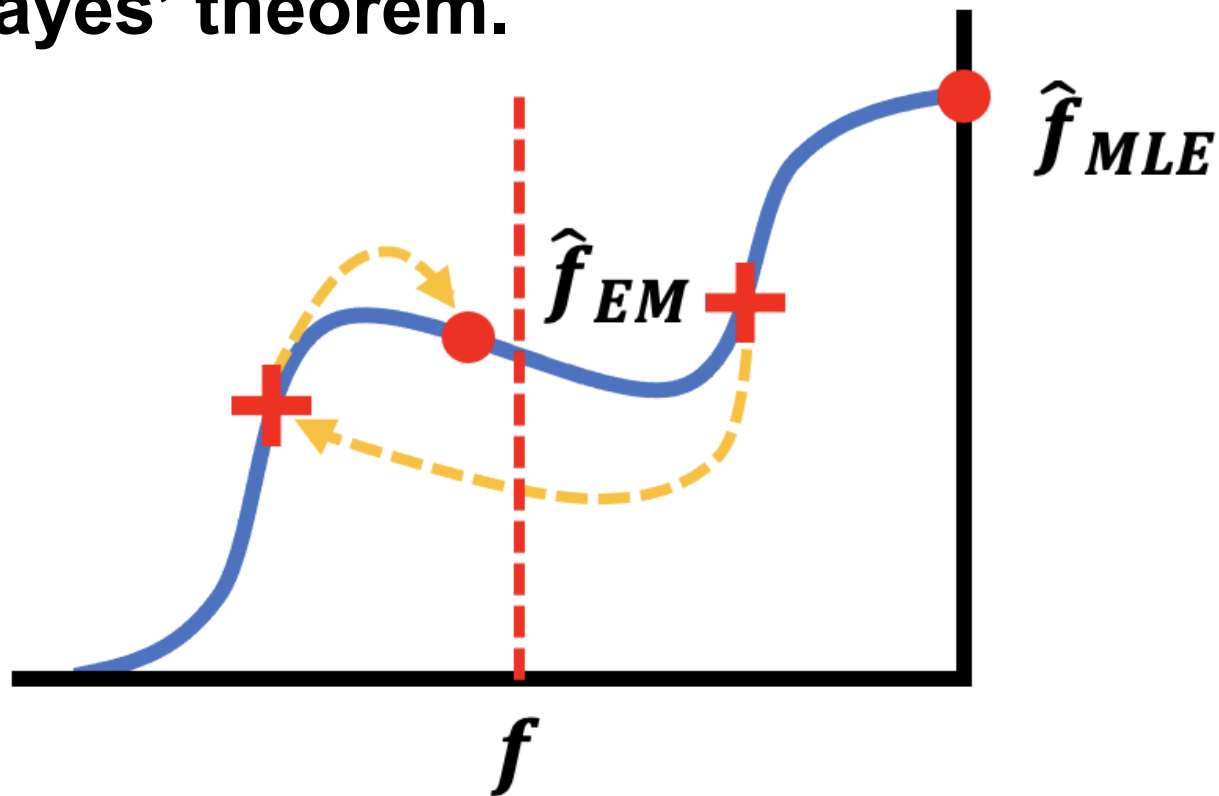
Our idea

- Our Approach
 - Expectation Maximization (EM) algorithm estimation
- Proposal
 - We apply the EM algorithm to the randomization used in PrivKV

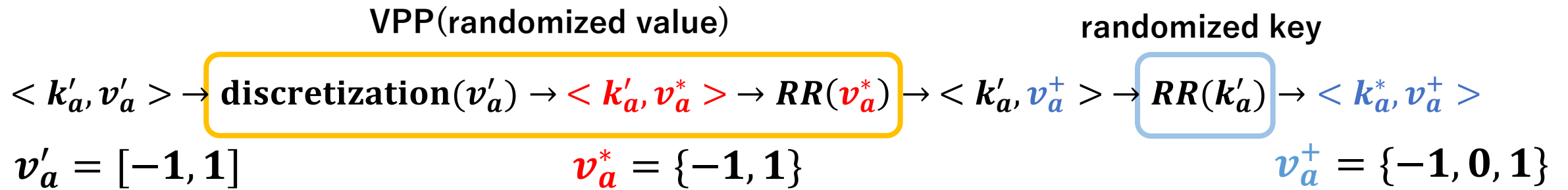
	Randomization	Estimation
PrivKV(M)	Randomized Response (RR)	MLE
Ours	Value Perturbation Primitive (VPP) +	EM

Why EM algorithm works better?

- EM algorithm [Dempster, et al., 1997]
- Iterative process for which posterior probabilities are updated based on Bayes' theorem.



Our idea



Estimating the marginal probability of $\langle k'_a, v_a^* \rangle$ using the output $\langle k_a^*, v_a^+ \rangle$

✗ $v_a^+ \rightarrow v'_a$ ✓ $v_a^+ \rightarrow v_a^*$

$$\langle k_a^*, v_a^+ \rangle = \{ \langle 1, 1 \rangle, \langle 1, -1 \rangle, \langle 0, 0 \rangle \}$$

$$\langle k'_a, v_a^* \rangle = \{ \langle 1, 1 \rangle, \langle 1, -1 \rangle, \langle 0, 1 \rangle, \langle 0, -1 \rangle \}$$

$$\langle k_a, v_a \rangle = \langle 0, _ * \rangle$$

Instead of independently calculating the frequency and mean values, we consider the key-value data as a whole

Experiment 1

- **Datasets**

 - 3 synthetic datasets

Distribution	#users	#items	value range
Gaussian	10^5	100	-1 - 1
Power-law			
Linear			

 - 2 open datasets

Datasets	#ratings	#users	#items	value range
MovieLens	10,000,054	69,877	10,677	0.5 - 5
Clothing	192,544	96,57	3,831	1 - 10

- **Evaluation metrics**

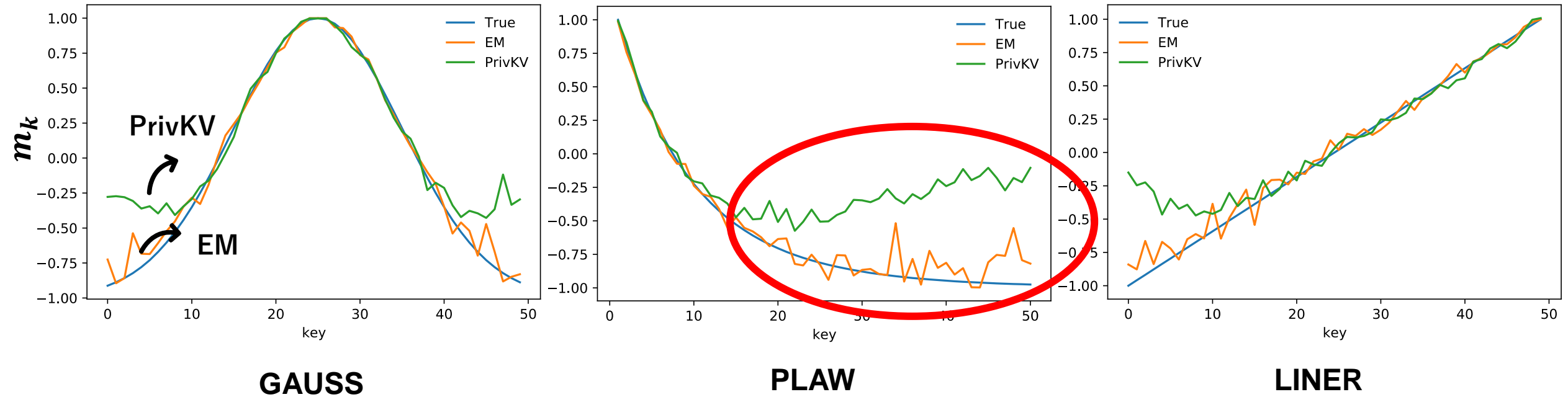
 - **Mean Square Error(MSE)**

$$MSE(f) = \frac{1}{|K|} \sum_{i=1}^{|K|} \left(\frac{\hat{f}_i}{n} - \frac{f_i}{n} \right)^2$$

$$MSE(m) = \frac{1}{|K|} \sum_{i=1}^{|K|} (\hat{m}_i - m_i)^2$$

Result1 Estimated mean distribution

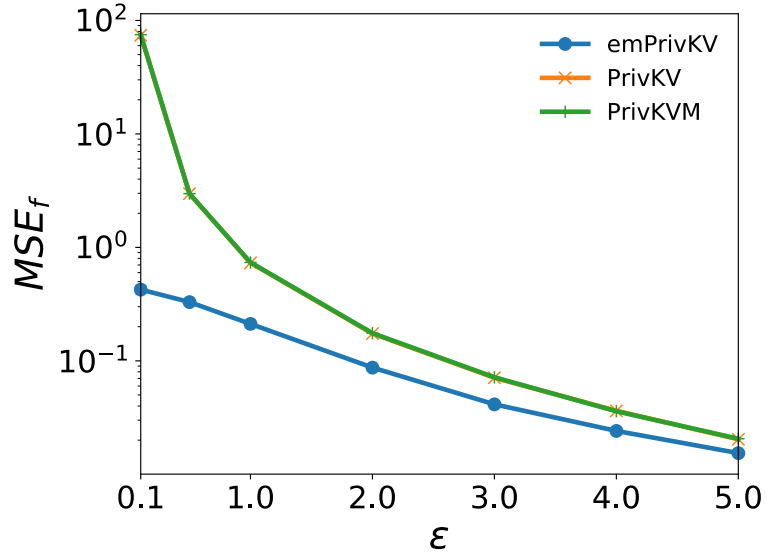
($\epsilon = 4, n = 10^5$)



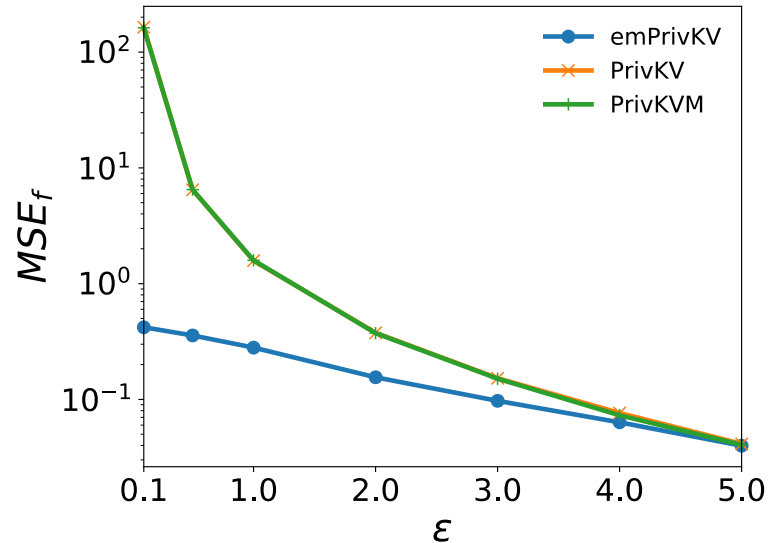
EM is more accurately when frequencies are small

Result 2 Estimation error MSE

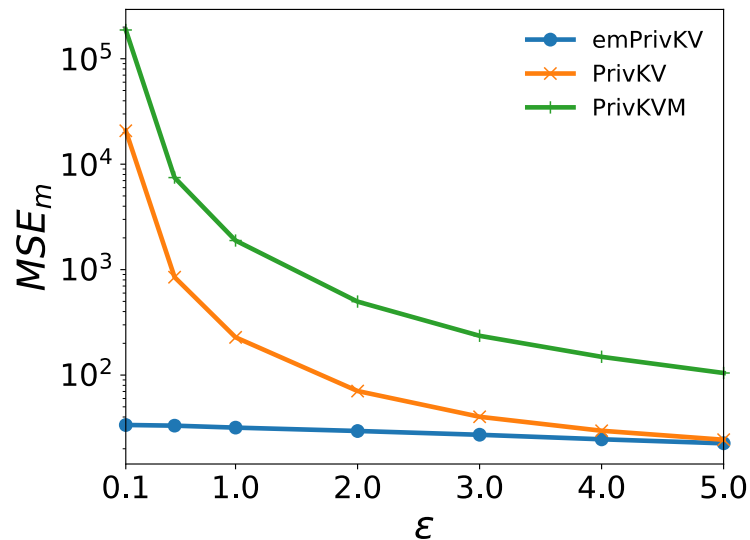
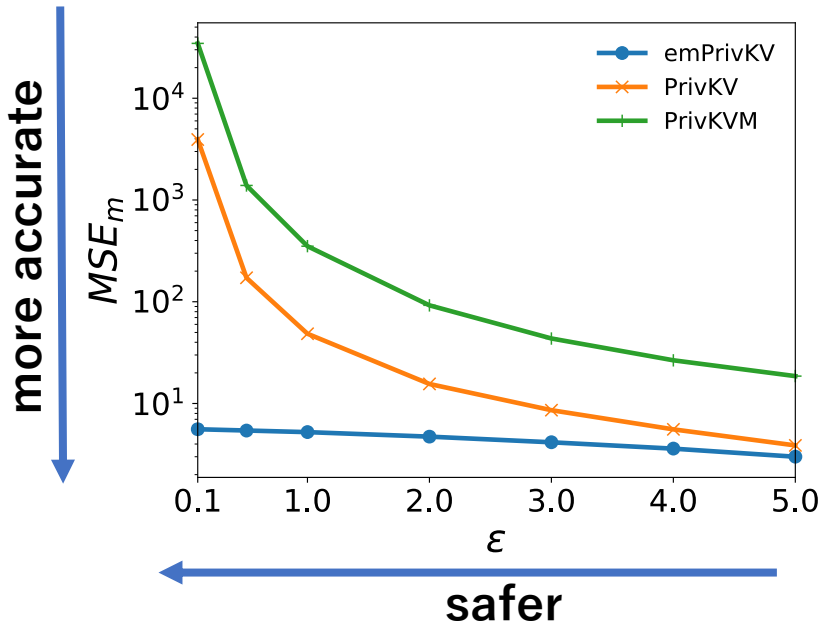
MovieLens



Clothing



- A proposed method is more accurate than the other two methods regardless to ϵ .

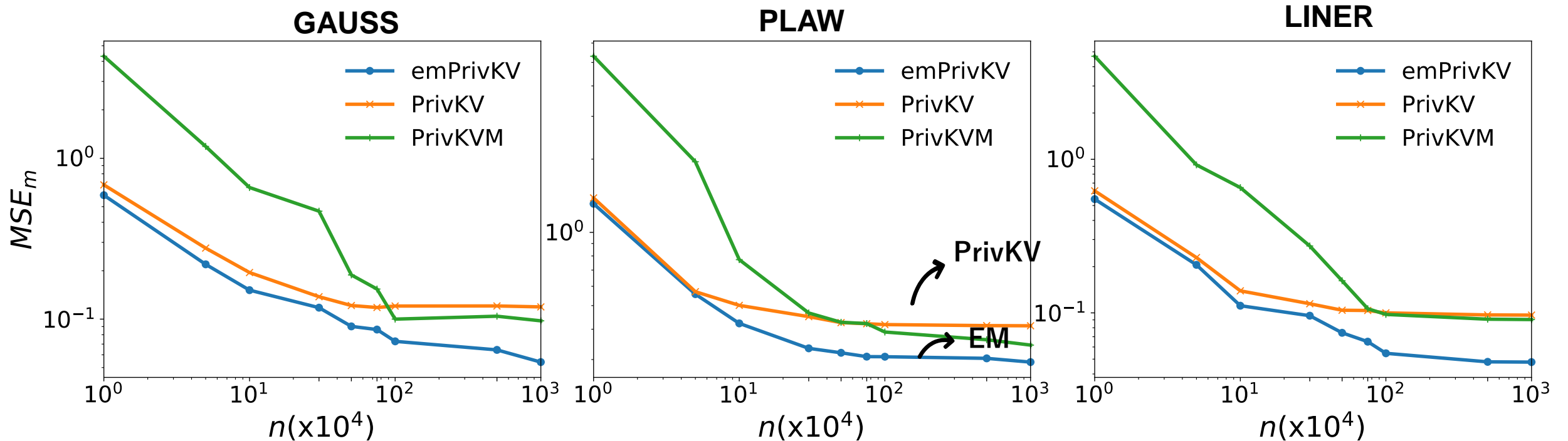


- Especially effective for smaller values of ϵ .

Result 3

MSE_m with regard to the number of users n

#key=50, $\epsilon = 2$



EM's MSE_m is the smallest for all n

Conclusion

- We studied the LDP algorithms for key-value data that estimate the key frequencies and the mean values.
- We propose an algorithm based on the EM algorithms to improve the estimation accuracy perturbed in the LDP algorithms PrivKV.
- The EM algorithm estimates the frequency and mean well even when the frequency of the key-value record is limited.
- Hence, we conclude that the proposed EM algorithm is appropriate for private data analysis dealing with rare item.