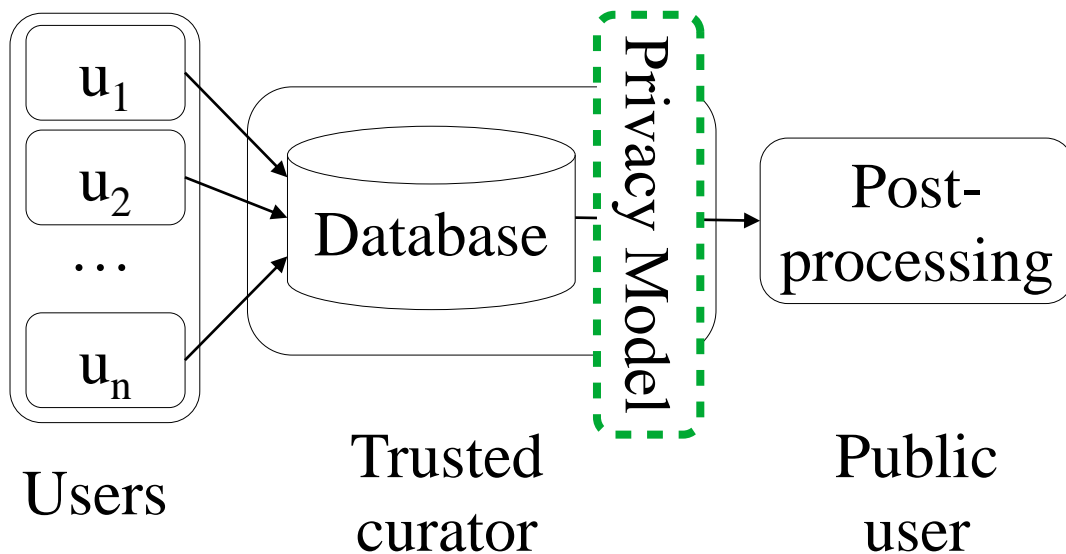


Random Sampling Plus Fake Data (RS+FD): Multidimensional Frequency Estimates With Local Differential Privacy*

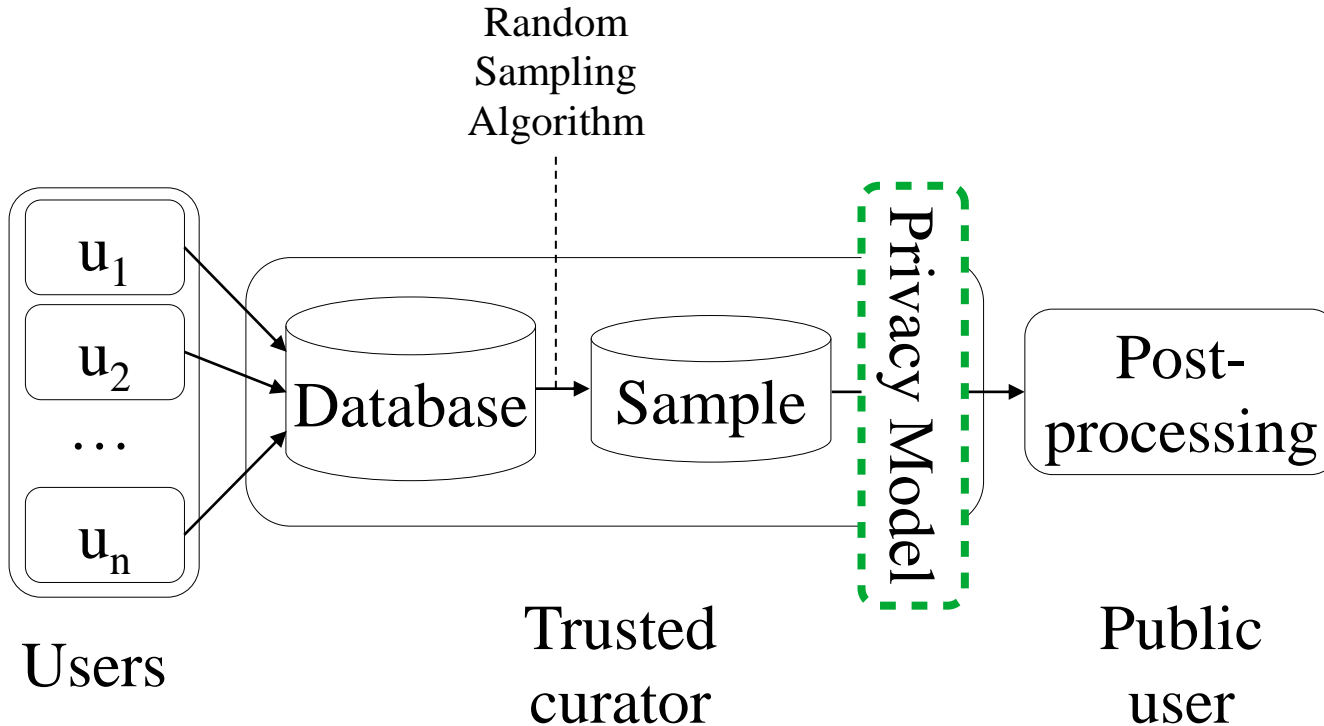
Héber H. ARCOLEZI, Jean-François COUCHOT, Bechara AL BOUNA, Xiaokui XIAO



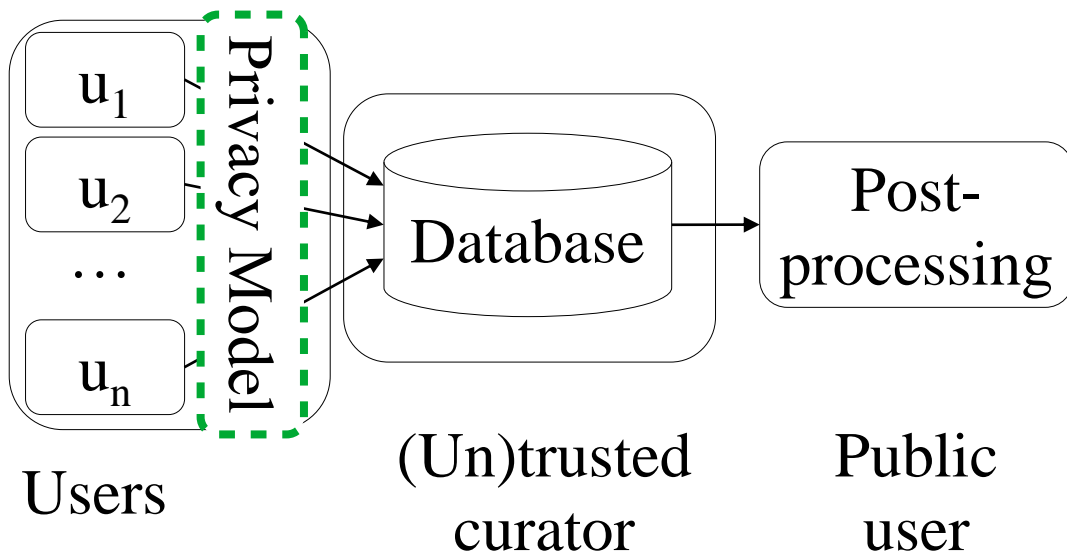


- *Centralized* setting of DP.
- Interpretation: The addition (or removal) of anyone's record has a minimal (ϵ) influence on the outcome.
- Small $\epsilon \rightarrow$ stronger privacy
- $\epsilon \rightarrow$ a.k.a. "privacy budget"
- Robust to post-processing.

Amplification by Sampling²



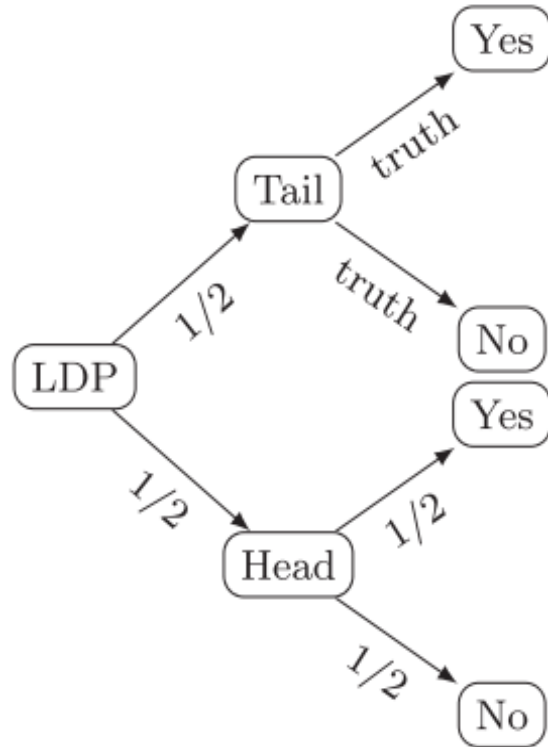
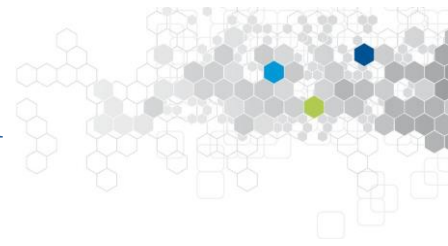
- Randomly subsample the database w/ sampling rate β .
- Interpretation: an attacker is unable to distinguish which data samples were used in the analysis.
- Amplification: $\epsilon' \geq \epsilon$
- $\epsilon = \ln(1 + \beta(e^{\epsilon'} - 1))$



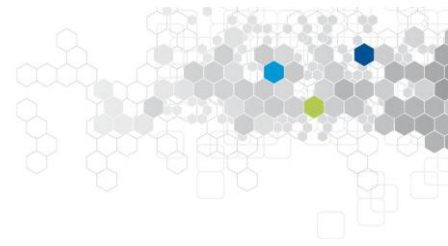
- *Local* setting of DP.
- Interpretation: Any two items have close probability (controlled by ϵ) to be mapped to the same perturbed value.
- Several LDP implementations in practice.



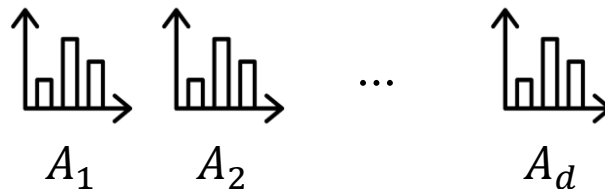
- Motivated by surveying people on sensitive topics.
- Main idea → Providing deniability to users' answer (yes/no → binary).
- Survey people: “Are you a member of the communist party?”
- Each person:
 - Throw a secret coin:
 - If tail throw the coin again (ignoring the outcome) and answer the question honestly.
 - If head, then throw the coin again and answer “Yes” if head, “No” if tail.



- $O_y \rightarrow$ proportion of *observed* yes
- $O_y \approx \frac{1}{2}t_y + \frac{1}{4}n$
- $t_y \rightarrow$ proportion of *true* yes
- $t_y \approx 2O_y - \frac{1}{2}n$
- Satisfies LDP w/:
 - $\epsilon = \ln(0.75/0.25) = \ln(3)$
 - prob. of 'being honest'
 - prob. of 'lying'



- **Key Issue:** Collecting *multidimensional* data under ϵ -LDP for the fundamental task of *frequency estimation*.
- **More formally (notation):**
 - d attributes $A = \{A_1, A_2, \dots, A_d\}$;
 - Each attribute A_j has a discrete domain D_j of size $|D_j| = k_j$;
 - Each user u_i for $1 \leq i \leq n$ has a tuple $v^i = (v_1^i, v_2^i, \dots, v_d^i)$;
 - **Analyzer:** estimate a k_j -bins histogram for each attribute $j \in [1, d]$.



- **Generalized RR (GRR)⁵**: Extends RR to the case of $k_j \geq 2$.

$$\forall_{y \in D_j} Pr[\psi_{GRR(\epsilon)}(v) = y] = \begin{cases} p = \frac{e^\epsilon}{e^\epsilon + k_j - 1}, & \text{if } y = v \\ q = \frac{1}{e^\epsilon + k_j - 1}, & \text{if } y \neq v \end{cases} \quad \epsilon = \ln\left(\frac{p}{q}\right)$$

- **Optimized Unary Encoding (OUE)⁶**: Encode as a bit-vector B and perturb each bit independently into a new bit-vector B' . More specifically:

$$Pr[B'_i = 1] = \begin{cases} p = 1/2, & \text{if } B_i = 1 \\ q = \frac{1}{e^\epsilon + 1}, & \text{if } B_i = 0 \end{cases} \quad \epsilon = \ln\left(\frac{p(1-q)}{q(1-p)}\right)$$

⁵ Kairouz, P., Bonawitz, K. and Ramage, D. Discrete distribution estimation under local privacy. In International Conference on Machine Learning (2016).

⁶ Wang, T., Blocki, J., Li, N. and Jha, S. Locally differentially private protocols for frequency estimation. In 26th USENIX Security Symposium (2017).

- **Unbiased Estimator:** To estimate the frequency $f(v_i)$ that a value v_i occurs for $i \in [1, k_j]$ one calculates

$$\hat{f}(v_i) = \frac{N_i - nq}{n(p-q)}, \quad N_i = \text{number of times the value (or bit) } i \text{ has been reported.}$$

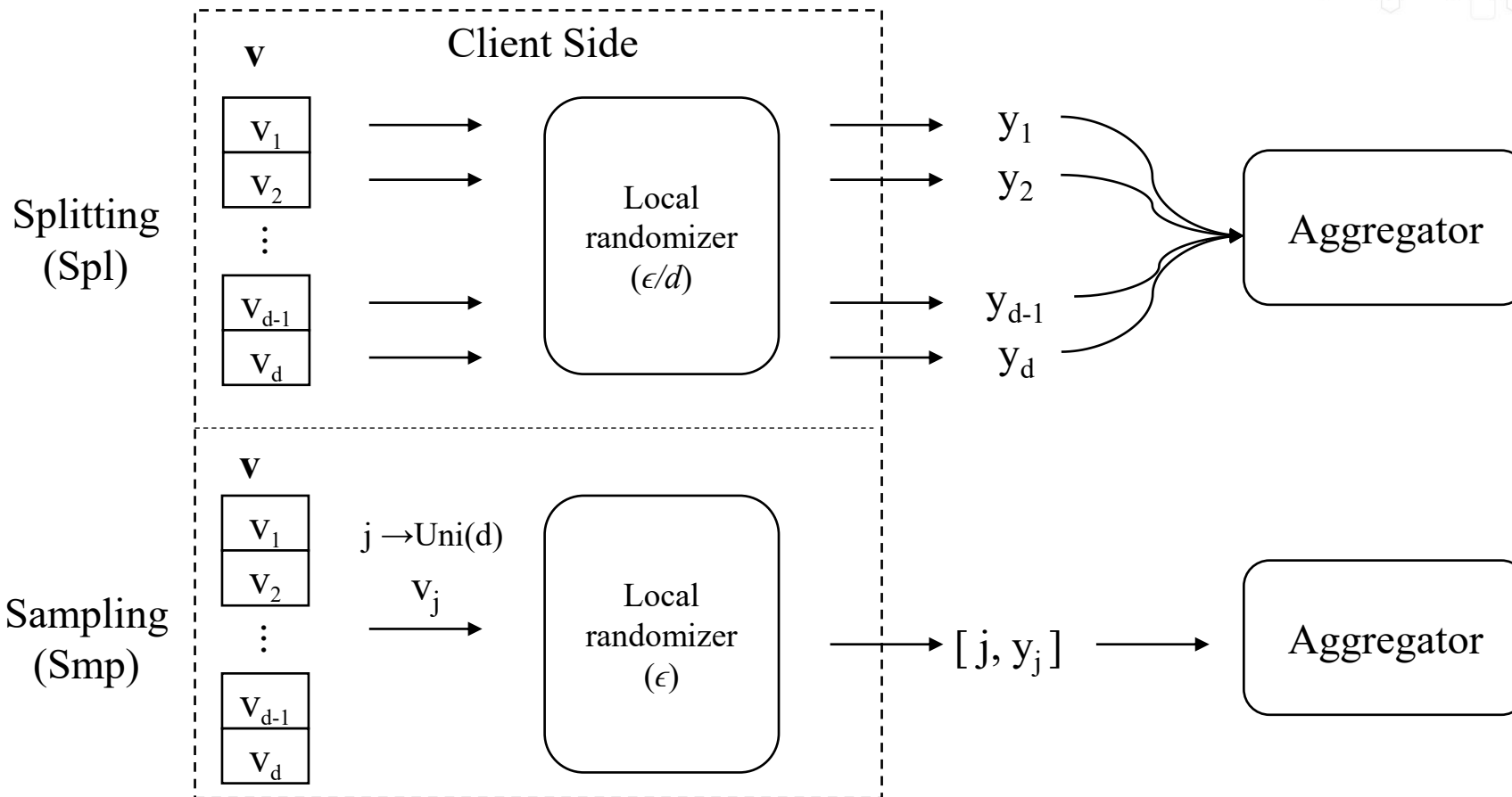
- **Approximate Variances:**

$$\text{Var}[\hat{f}_{GRR}(v_i)] = \frac{e^\epsilon + k_j - 2}{n(p-q)^2} \qquad \text{Var}[\hat{f}_{OUE}(v_i)] = \frac{4e^\epsilon}{n(p-q)^2}$$

- **Adaptive LDP protocol⁶:** Given $k_j, p, q,$ and ϵ

$$ADP = \begin{cases} GRR & \text{if } k_j < 3e^\epsilon + 2 \\ OUE & \text{otherwise.} \end{cases}$$

What is the state-of-the-art?



Why not *Smp*?

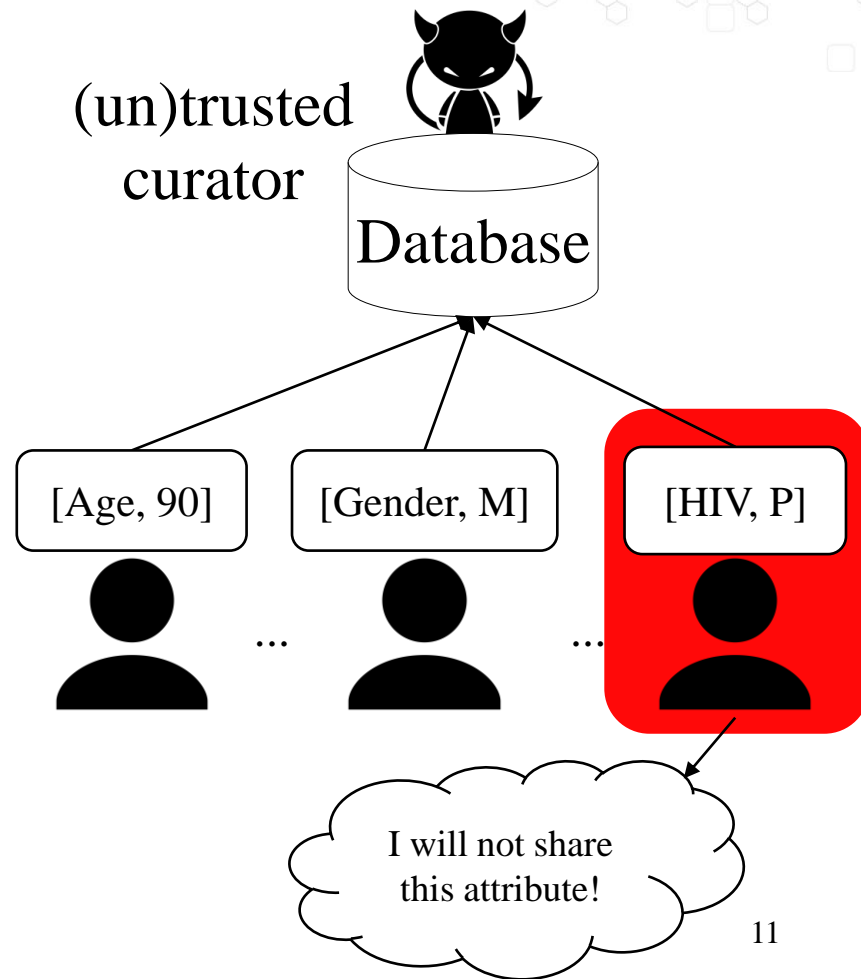
Example:

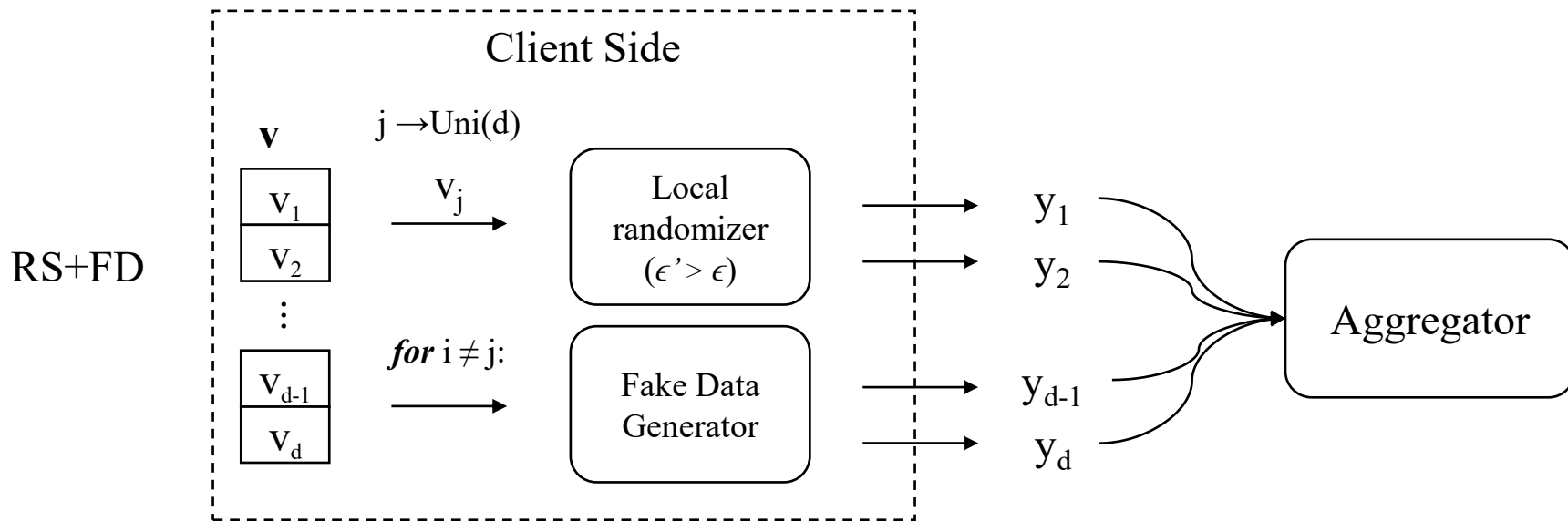
GRR for attributes with small domain
OUE otherwise

- $Smp[ADP] \rightarrow (\text{attribute}, \epsilon\text{-LDP value})$
- Application scenario: health data
- $\epsilon = 2$, $d = 3$ attributes: age ($k_1 = [1, \dots, 100]$), gender ($k_2 = [M, F]$), and HIV ($k_3 = [P, N]$).

$$p_{grr} = \frac{e^\epsilon}{e^\epsilon + k_j - 1} \approx 0.88 \text{ (probability of 'being honest')}$$

$$q_{grr} = \frac{1 - p_{grr}}{k_j - 1} \approx 0.12 \text{ (probability of 'lying')}$$





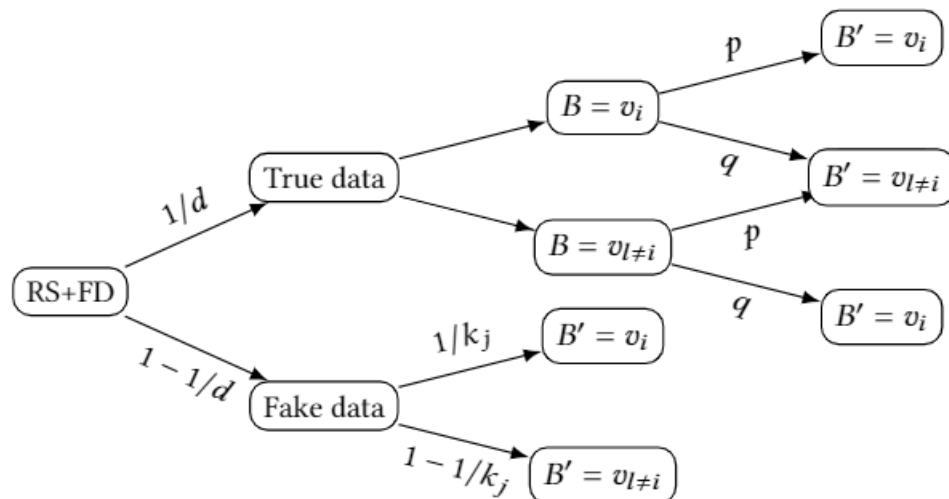


Algorithm 1 RS+FD[GRR]: Client Side

Input : tuple $\mathbf{v} = (v_1, v_2, \dots, v_d)$, domain size of attributes $\mathbf{k} = [k_1, k_2, \dots, k_d]$, privacy parameter ϵ , local randomizer GRR.

Output : privatized tuple $\mathbf{y} = (y_1, y_2, \dots, y_d)$.

- 1: $\epsilon' = \ln(d \cdot (e^\epsilon - 1) + 1)$ ▶ amplification by sampling [31]
 - 2: $j \leftarrow \text{Uniform}(\{1, 2, \dots, d\})$ ▶ Selection of attribute to privatize
 - 3: $B_j \leftarrow v_j$
 - 4: $y_j \leftarrow \text{GRR}(B_j, k_j, \epsilon')$ ▶ privatize data of the sampled attribute
 - 5: for $i \in \{1, 2, \dots, d\}/j$ do ▶ non-sampled attributes
 - 6: $y_i \leftarrow \text{Uniform}(\{1, \dots, k_i\})$ ▶ generate fake data
 - 7: end for
- ▶ sampling result is not disclosed



Aggregator \rightarrow For each attribute, estimate: $\hat{f}(v_i) = \frac{N_i d k_j - n(d - 1 + q k_j)}{n k_j (p - q)}$

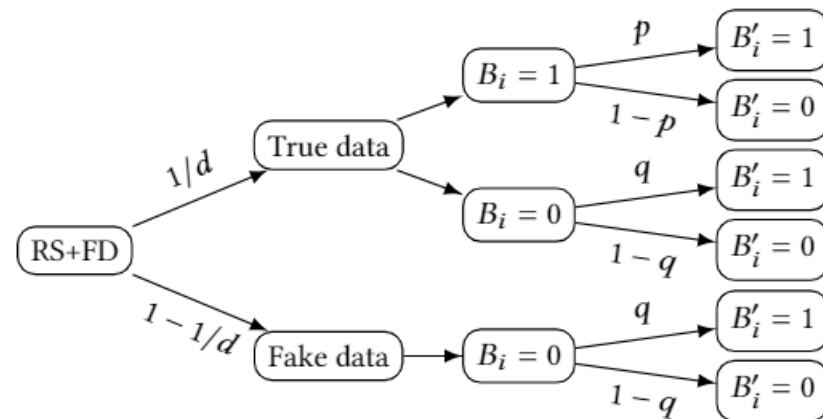


Algorithm 2 RS+FD[OUE-z]: Client Side

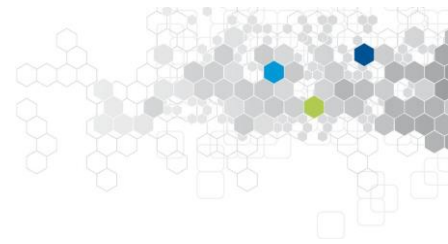
Input : tuple $\mathbf{v} = (v_1, v_2, \dots, v_d)$, domain size of attributes $\mathbf{k} = [k_1, k_2, \dots, k_d]$, privacy parameter ϵ , local randomizer OUE.

Output : privatized tuple $\mathbf{B}' = (B'_1, B'_2, \dots, B'_d)$.

- 1: $\epsilon' = \ln(d \cdot (e^\epsilon - 1) + 1)$
 - 2: $j \leftarrow \text{Uniform}(\{1, 2, \dots, d\})$
 - 3: $B_j = \text{Encode}(v_j) = [0, 0, \dots, 1, 0, \dots, 0]$
 - 4: $B'_j \leftarrow \text{OUE}(B_j, \epsilon')$
 - 5: **for** $i \in \{1, 2, \dots, d\} / j$ **do**
 - 6: $B_i \leftarrow [0, 0, \dots, 0]$
 - 7: $B'_i \leftarrow \text{OUE}(B_i, \epsilon')$
 - 8: **end for**
- return :** $\mathbf{B}' = (B'_1, B'_2, \dots, B'_d)$
- amplification by sampling [31]
 - Selection of attribute to privatize
 - one-hot-encoding
 - privatize real data with OUE
 - non-sampled attributes
 - initialize zero-vectors
 - randomize zero-vector with OUE
 - sampling result is not disclosed



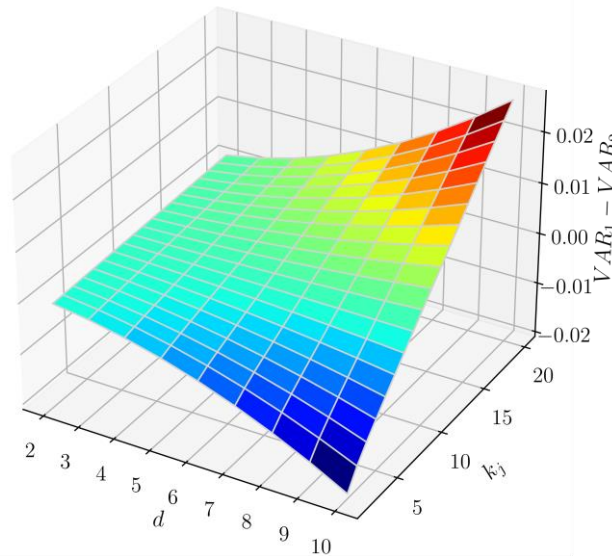
Aggregator \rightarrow For each attribute, estimate: $\hat{f}(v_i) = \frac{d(N_i - nq)}{n(p - q)}$

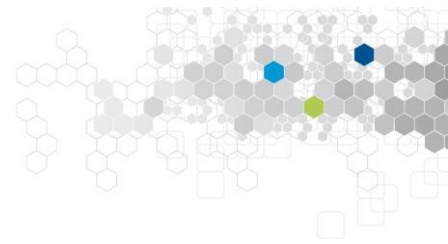


- Let $VAR_1 = VAR_{RS+FD}[GRR]$ and $VAR_2 = VAR_{RS+FD}[OUE-z]$
- For each attribute, given d , k_j , and ϵ' , select RS+FD[GRR] **if**:

$$VAR_1 \leq VAR_2, \text{ i.e., if } VAR_1 - VAR_2 \leq 0$$

- Let $n = 10000$, $d \in [2, 10]$, $k_j \in [2, 20]$, and $\epsilon' = \ln(3)$



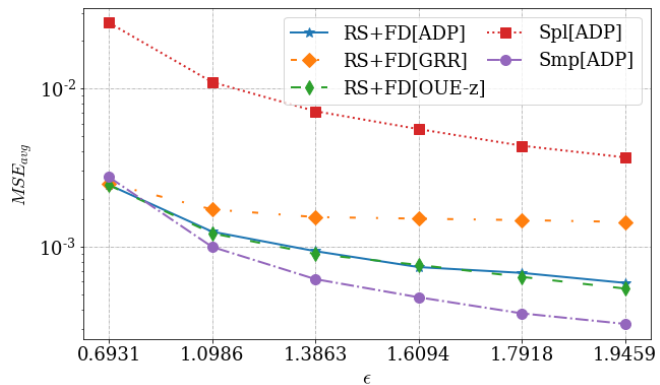
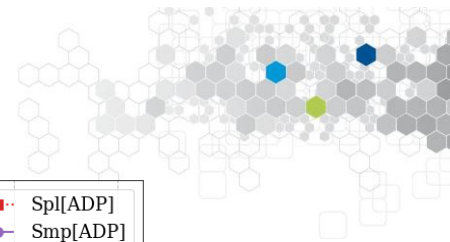


- Datasets:
 - Nursery⁷: $n = 12960$, $d = 9$, $\mathbf{k} = [3,5,4,4,3,2,3,3,5]$
 - Adults⁷: $n = 45422$, $d = 9$, $\mathbf{k} = [7,16,7,14,6,5,2,41,2]$
 - MS-FIMU⁸: $n = 88935$, $d = 6$, $\mathbf{k} = [3,3,8,12,37,11]$
 - Census-Income⁷: $n = 299285$, $d = 33$, $\mathbf{k} = [9,52,47,17, \dots, 3,3,2]$
- Evaluation: $\epsilon = [\ln(2), \ln(3), \dots, \ln(7)]$.
- Metric:

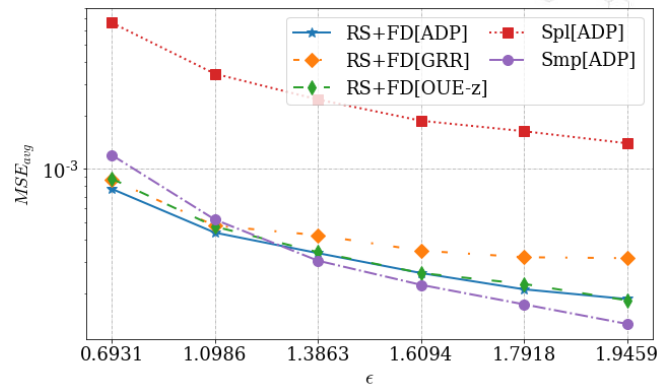
$$\text{MSE}_{avg} = \frac{1}{d} \sum_{j \in [1,d]} \frac{1}{|D_j|} \sum_{v_i \in D_j} (f(v_i) - \hat{f}(v_i))^2$$

⁷ Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository: <http://archive.ics.uci.edu/ml/index.php>

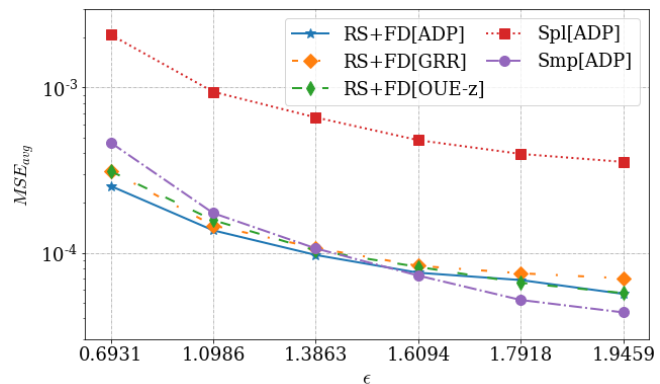
⁸ Arcolezi, H.H., Couchot, J.F., Baala, O., et al. Mobility modeling through mobile data: generating an optimized and open dataset respecting privacy. In 16th IWCMC (2020).



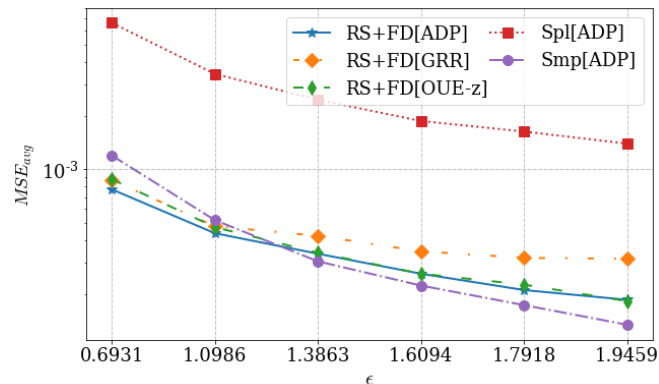
Nursery⁷: $n = 12960$, $d = 9$, $\mathbf{k} = [3,5,4,4,3,2,3,3,5]$



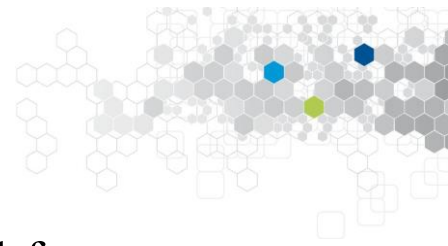
Adults⁷: $n=45422$, $d=9$, $\mathbf{k}=[7,16,7,14,6,5,2,41,2]$



MS-FIMU⁸: $n = 88935$, $d = 6$, $\mathbf{k} = [3,3,8,12,37,11]$



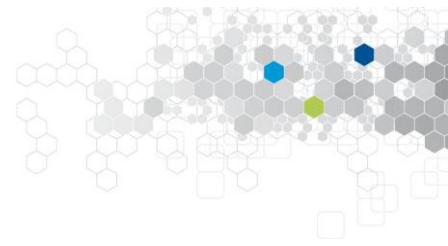
Census-Income⁷: $n = 299285$, $d = 33$, $\mathbf{k} = [9,52,47, \dots, 3,2]$



- We propose a generic framework RS+FD for multidimensional frequency estimates under LDP with theoretical proofs.
- RS+FD achieves nearly the same or better utility than *Smp* with higher privacy protection (uncertainty).
- Limitations:
 - Sampling error + noise from fake reports;
 - More computation and communication cost than *Smp*.
- Perspectives:
 - Cast other LDP protocols into RS+FD;
 - Attack: is it possible to state which attribute value is “fake”?



1. Dwork, C., Roth, A. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science* (3–4), 211–407 (2014).
2. Ninghui Li, Wahbeh Qardaji, and Dong Su. On sampling, anonymization, and differential privacy or, k-anonymization meets differential privacy. In *Proceedings of the 7th ACM Symposium on Information, Computer and Communications Security - ASIACCS'12* (2012).
3. Kasiviswanathan, S.P., Lee, H.K., Nissim, K., Raskhodnikova, S., Smith, A. What can we learn privately? In: *49th Annual IEEE Symposium on Foundations of Computer Science* (2008).
4. Warner, S.L. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 1965.
5. Kairouz, P., Bonawitz, K. and Ramage, D. Discrete distribution estimation under local privacy. In *International Conference on Machine Learning* (2016).
6. Wang, T., Blocki, J., Li, N. and Jha, S. Locally differentially private protocols for frequency estimation. In *26th USENIX Security Symposium* (2017).
7. Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository: <http://archive.ics.uci.edu/ml/index.php>
8. Arcolezi, H.H., Couchot, J.F., Baala, O., et al. Mobility modeling through mobile data: generating an optimized and open dataset respecting privacy. In *2020 International Wireless Communications and Mobile Computing* (2020).



Thank you very much for your attention!!!

Questions?

Codes → <https://github.com/hharcolezi/ldp-protocols-mobility-cdrs>

Contact → heber.hwang_arcolezi@univ-fcomte.fr