

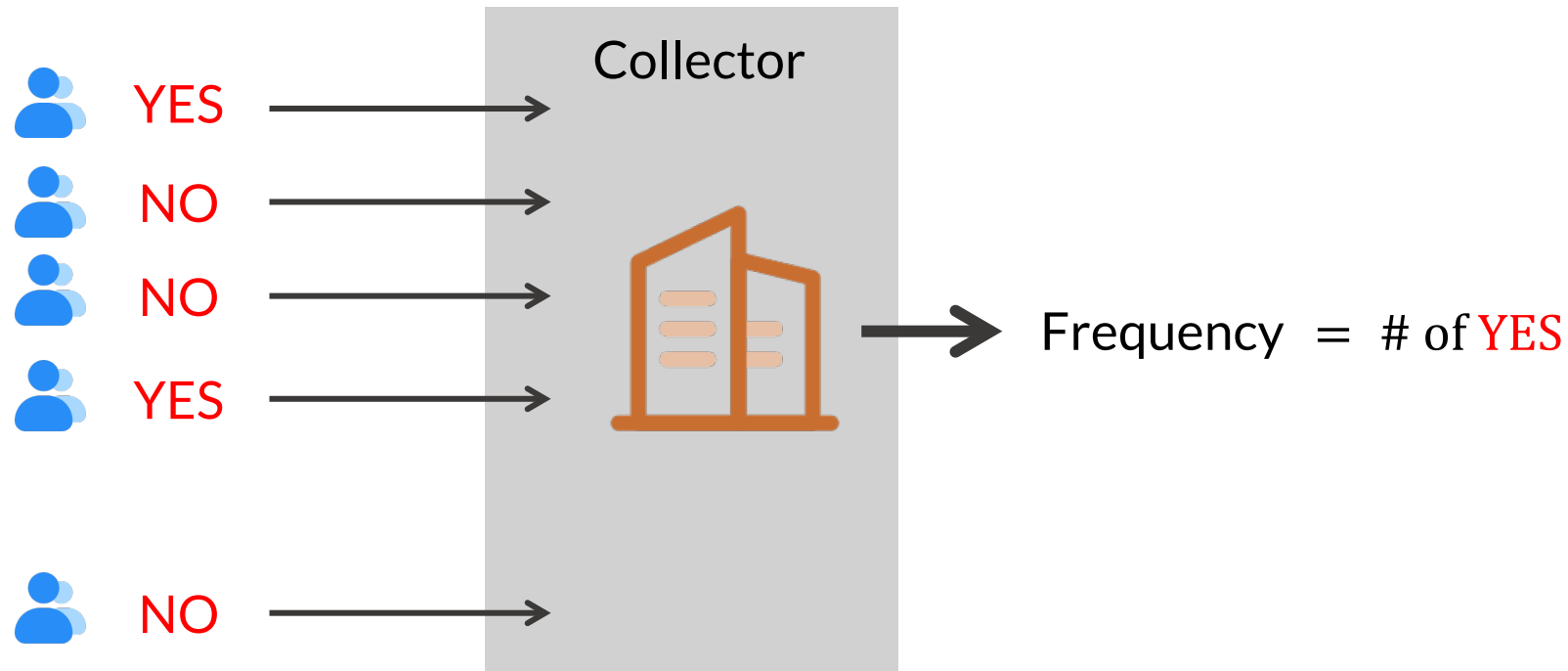
# Locally Differentially Private Frequency Estimation via Joint Randomized Response

Authors: [Ye Zheng](#), Shafizur Rahman Seeam, Yidan Hu, [Rui Zhang](#), [Yanchao Zhang](#)



# Frequency Estimation

- Social scientists: How many people engage in tax evasion?
  - ask one person if they had evaded tax
  - the person answers **YES** or **NO**



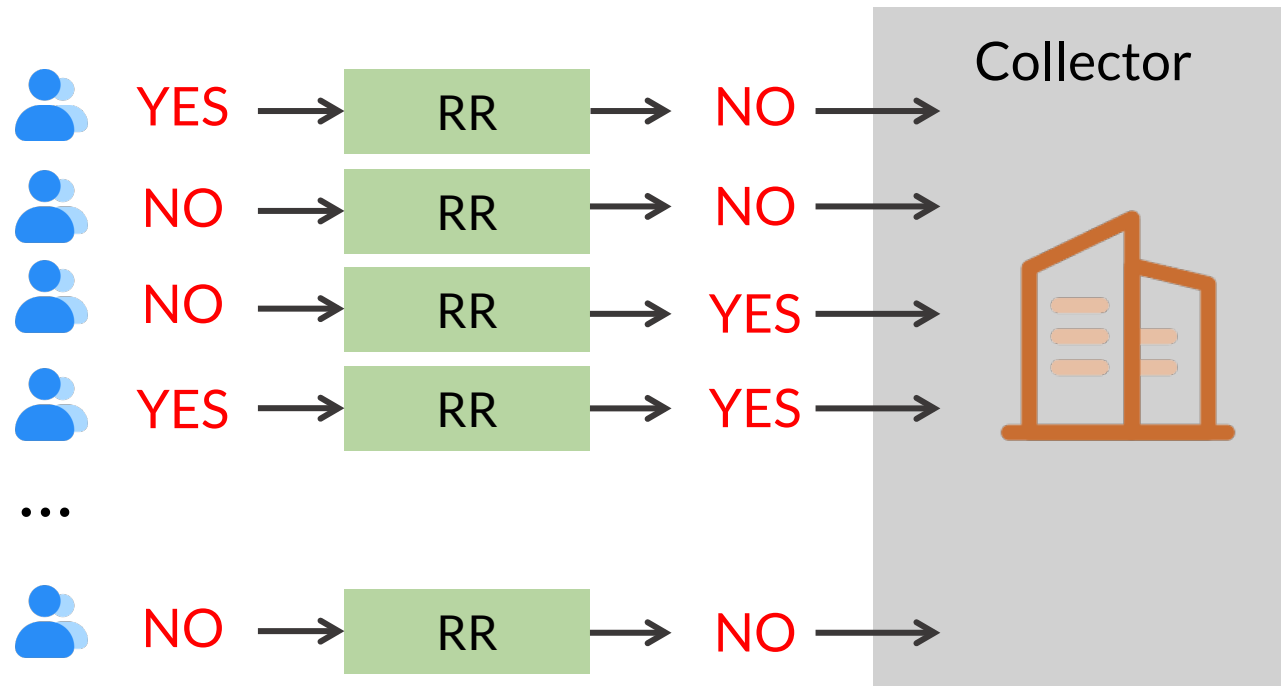
# Randomized Response for Privacy

---

- People **tend to lie** for such sensitive/embarrassing question
  - i.e. don't want to let the collector know

# Randomized Response for Privacy

- People **tend to lie** for such sensitive/embarrassing question  
- i.e. don't want to let the collector know
- Randomized Response: Randomize the truth **before answering the collector**

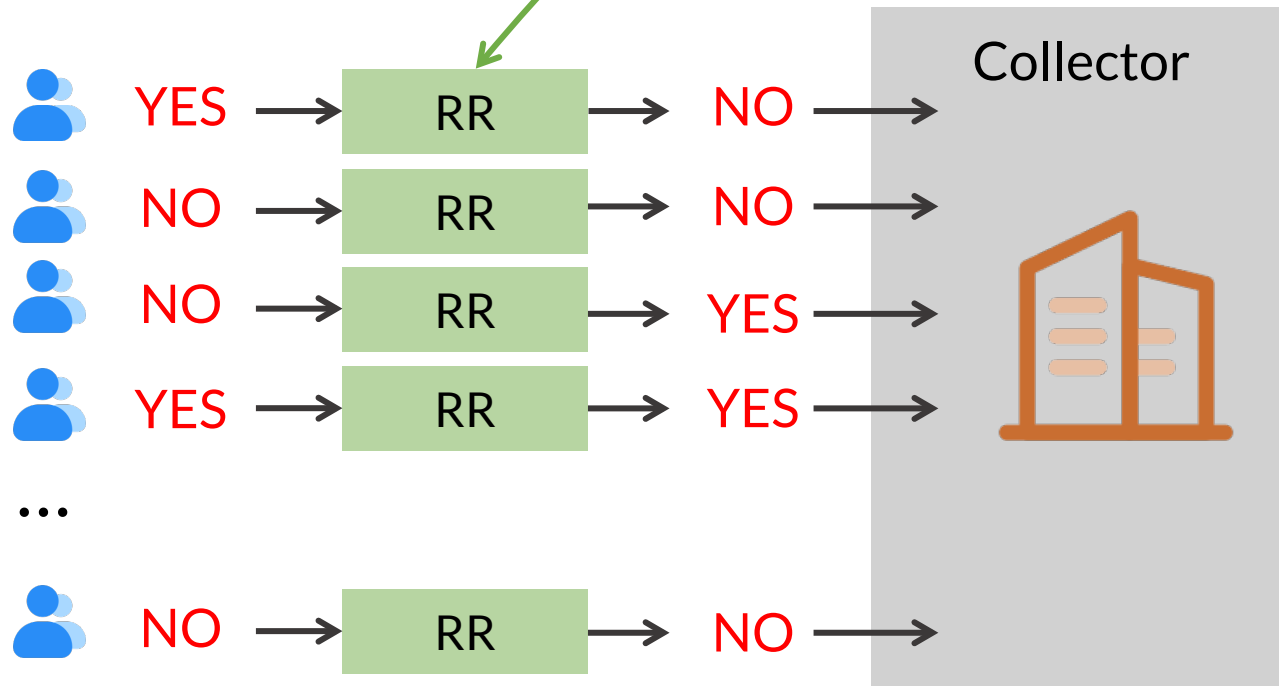


# Randomized Response for Privacy

- People **tend to lie** for such sensitive/embarrassing question  
- i.e. don't want to let the collector know
- Randomized Response: Randomize the truth **before answering**

**RR:**  
answer truth with probability  $p$

$$RR(x) = \begin{cases} x & \text{w.p. } p \\ \neg x & \text{w.p. } 1 - p \end{cases}$$

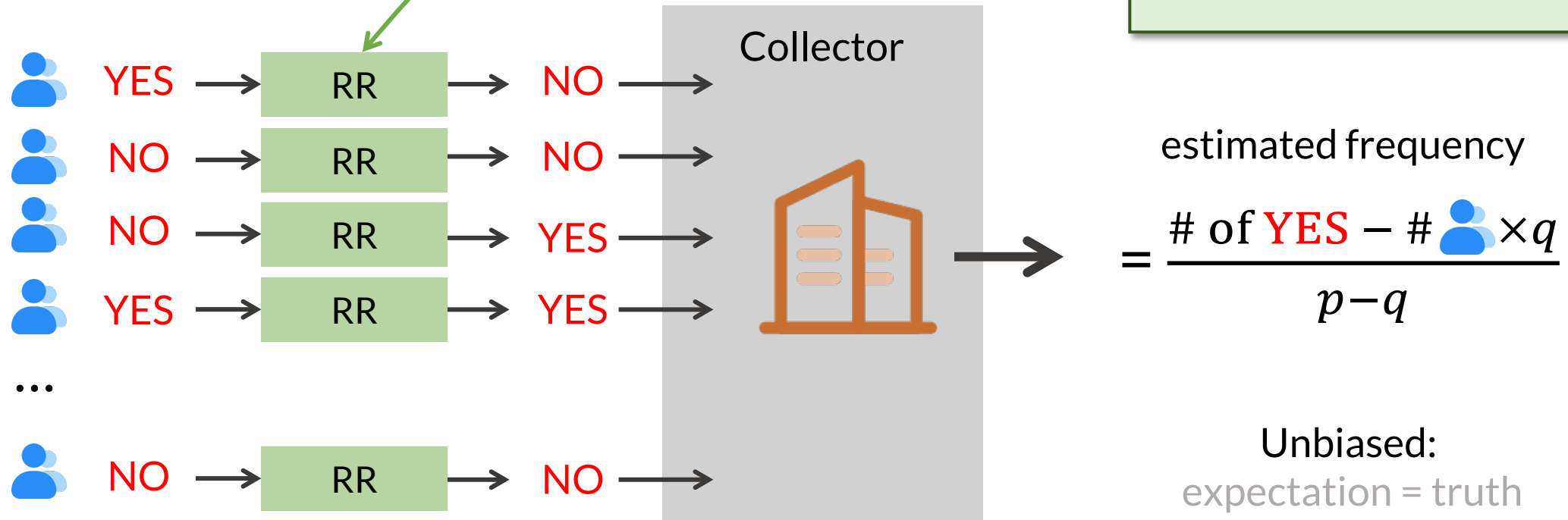


# Randomized Response for Privacy

- People **tend to lie** for such sensitive/embarrassing question  
- i.e. don't want to let the collector know
- Randomized Response: Randomize the truth **before answering**

**RR:**  
answer truth with probability  $p$

$$RR(x) = \begin{cases} x & \text{w.p. } p \\ \neg x & \text{w.p. } 1 - p \end{cases}$$



# Privacy: RR Satisfies LDP

- A mechanism  $\mathcal{M}$  satisfies LDP if

For any truth  $x_1, x_2$ ,  
and randomized answer  $y$ :

$$\max \frac{\Pr[\mathcal{M}(x_1) = y]}{\Pr[\mathcal{M}(x_2) = y]} \leq e^\epsilon$$

Distinguishability of  $x_1$  (YES) and  $x_2$  (NO)  
from  $y$  (randomized answer)


# Privacy: RR Satisfies LDP

- A mechanism  $\mathcal{M}$  satisfies LDP if

For any truth  $x_1, x_2$ ,  
and randomized answer  $y$ :

$$\max \frac{\Pr[\mathcal{M}(x_1) = y]}{\Pr[\mathcal{M}(x_2) = y]} \leq e^\epsilon$$

Distinguishability of  $x_1$  (YES) and  $x_2$  (NO)  
from  $y$  (randomized answer)


$$\text{RR: } \frac{p}{q} \rightarrow \epsilon \geq \ln \frac{p}{q}$$



# Privacy: RR Satisfies LDP



- A mechanism  $\mathcal{M}$  satisfies LDP if

For any truth  $x_1, x_2$ ,  
and randomized answer  $y$ :

$$\max \frac{\Pr[\mathcal{M}(x_1) = y]}{\Pr[\mathcal{M}(x_2) = y]} \leq e^\epsilon$$

Distinguishability of  $x_1$  (YES) and  $x_2$  (NO)  
from  $y$  (randomized answer)

RR:  $\frac{p}{q} \rightarrow \epsilon \geq \ln \frac{p}{q}$

- **quantifiable hardness** to distinguish  $x_1$  (YES) and  $x_2$  (NO) from the randomized answer  $y$
- against inference from data collectors  or adversaries 

# Utility: RR's Variance

- Randomization reduces data utility

$$\text{Var}\left[\frac{\# \text{ of YES} - \# \text{ of people} \times q}{p - q}\right] = \frac{\text{Var}[\# \text{ of YES}]}{(p - q)^2} = \frac{npq}{(p - q)^2}$$

- summation of variance from  $n$  independent randomization

# Utility: RR's Variance

- Randomization reduces data utility

$$\text{Var}\left[\frac{\# \text{ of YES} - \# \text{ of people} \times q}{p - q}\right] = \frac{\text{Var}[\# \text{ of YES}]}{(p - q)^2} = \frac{npq}{(p - q)^2}$$

- summation of variance from  $n$  independent randomization

- larger  $p \in (0.5, 1]$   $\rightarrow$  lower variance  $\rightarrow$  larger privacy parameter  $\epsilon$

$\uparrow$  data utility



$\downarrow$  privacy

# Utility: RR's Variance

- Randomization reduces data utility

$$\text{Var}\left[\frac{\# \text{ of YES} - \# \text{ of people} \times q}{p - q}\right] = \frac{\text{Var}[\# \text{ of YES}]}{(p - q)^2} = \frac{npq}{(p - q)^2}$$

- summation of variance from  $n$  independent randomization

- larger  $p \in (0.5, 1]$   $\rightarrow$  lower variance  $\rightarrow$  larger privacy parameter  $\epsilon$

$\uparrow$  data utility



$\downarrow$  privacy

- Q: Can correlated (joint) randomization improve this privacy-utility tradeoff?

# This Paper: Joint RR (JRR)

---

- Joint randomization can boost data utility

# This Paper: Joint RR (JRR)

- Joint randomization can boost data utility
- Example:** 2-person ( $x_1 = \text{YES}$  and  $x_2 = \text{YES}$ ) with  $p = 0.8$  ( $P[T = 1] = 0.8$ )

RR: Joint distribution

	$T_1 = 1$	$T_1 = 0$	Truthfulness of $x_1$
$T_2 = 1$	0.64 ( $=0.8 \times 0.8$ )	0.16 ( $=0.2 \times 0.8$ )	
$T_2 = 0$	0.16 ( $=0.8 \times 0.2$ )	0.04 ( $=0.2 \times 0.2$ )	

Truthfulness  
of  $x_2$

# This Paper: Joint RR (JRR)

- Joint randomization can boost data utility
- **Example:** 2-person ( $x_1 = \text{YES}$  and  $x_2 = \text{YES}$ ) with  $p = 0.8$  ( $P[T = 1] = 0.8$ )

RR: Joint distribution

	$T_1 = 1$	$T_1 = 0$	Truthfulness of $x_1$
$T_2 = 1$	0.64 ( $=0.8 \times 0.8$ )	0.16 ( $=0.2 \times 0.8$ )	
$T_2 = 0$	0.16 ( $=0.8 \times 0.2$ )	0.04 ( $=0.2 \times 0.2$ )	

Truthfulness  
of  $x_2$

**Independent  $T_1$  and  $T_2$**  ( $P[T_1 \cap T_2] = P[T_1] \cdot P[T_2]$ )

# This Paper: Joint RR (JRR)

- Joint randomization can boost data utility
- Example:** 2-person ( $x_1 = \text{YES}$  and  $x_2 = \text{YES}$ ) with  $p = 0.8$  ( $P[T = 1] = 0.8$ )

RR: Joint distribution

	$T_1 = 1$	$T_1 = 0$
$T_2 = 1$	0.64 ( $=0.8 \times 0.8$ )	0.16 ( $=0.2 \times 0.8$ )
$T_2 = 0$	0.16 ( $=0.8 \times 0.2$ )	0.04 ( $=0.2 \times 0.2$ )

Truthfulness  
of  $x_2$

**Independent  $T_1$  and  $T_2$**  ( $P[T_1 \cap T_2] = P[T_1] \cdot P[T_2]$ )

JRR: Joint distribution

	$T_1 = 1$	$T_1 = 0$
$T_2 = 1$	0.6	0.2
$T_2 = 0$	0.2	0



# This Paper: Joint RR (JRR)

- Joint randomization can boost data utility
- Example:** 2-person ( $x_1 = \text{YES}$  and  $x_2 = \text{YES}$ ) with  $p = 0.8$  ( $P[T = 1] = 0.8$ )

RR: Joint distribution

	$T_1 = 1$	$T_1 = 0$
$T_2 = 1$	0.64 ( $=0.8 \times 0.8$ )	0.16 ( $=0.2 \times 0.8$ )
$T_2 = 0$	0.16 ( $=0.8 \times 0.2$ )	0.04 ( $=0.2 \times 0.2$ )

Truthfulness  
of  $x_2$

**Independent  $T_1$  and  $T_2$**  ( $P[T_1 \cap T_2] = P[T_1] \cdot P[T_2]$ )

JRR: Joint d

	$T_1 = 1$	$T_1 = 0$
$T_2 = 1$	0.6	0.2
$T_2 = 0$	0.2	0

$P[T_1 = 1] = 0.6 + 0.2 = 0.8$

$P[T_1 = 0 \cap T_2 = 0] = 0$   
 $P[T_1 = 0] \cdot P[T_2 = 0] = 0.04$

**NOT independent  $T_1$  and  $T_2$**

# Utility: JRR's Variance

- Same estimator as RR

$$\text{Expectation: } E[\# \text{ of YES}] = \sum_{i=1}^{\# \text{ people}} P[y_i = \text{YES}] = n_{\text{YES}} \cdot p + (\# \text{ people} - n_{\text{YES}}) \cdot q$$

$$\hat{n}_{\text{YES}} = \frac{\# \text{ of YES} - 2q}{p - q} \text{ is unbiased}$$

Identical to RR

# Utility: JRR's Variance

- Same estimator as RR

$$\text{Expectation: } E[\# \text{ of YES}] = \sum_{i=1}^{\# \text{ people}} P[y_i = \text{YES}] = n_{\text{YES}} \cdot p + (\# \text{ people} - n_{\text{YES}}) \cdot q$$

$$\hat{n}_{\text{YES}} = \frac{\# \text{ of YES} - 2q}{p - q} \text{ is unbiased}$$

Identical to RR

- Variance: ( $\# \text{ people} = 2, p = 0.8$ )

$$\text{Var}[\hat{n}_{\text{YES}}] = \frac{\text{Var}[\# \text{ of YES}]}{(0.8 - 0.2)^2}$$

JRR

# of YES	0	1	2
Probability	0	0.2 + 0.2	0.6

$$\text{Var}[\# \text{ of YES}] = E[(X - \mu)^2] = 0.24$$

# Utility: JRR's Variance

- Same estimator as RR

Expectation:  $E[\# \text{ of YES}] = \sum_{i=1}^{\# \text{ people}} P$

RR			
# of YES	0	1	2
Probability	0.04	0.16 + 0.16	0.6

$\text{Var}[\# \text{ of YES}] = E[(X - \mu)^2] \approx 0.32$

- Variance: ( $\# \text{ people} = 2, p = 0.8$ )

$$\text{Var}[\hat{n}_{\text{YES}}] = \frac{\text{Var}[\# \text{ of YES}]}{(0.8 - 0.2)^2}$$

JRR			
# of YES	0	1	2
Probability	0	0.2 + 0.2	0.6

$$\text{Var}[\# \text{ of YES}] = E[(X - \mu)^2] = 0.24$$

# Utility: JRR's Variance

- Same estimator as RR

Expectation:  $E[\# \text{ of YES}] = \sum_{i=1}^{\# \text{ people}} P$

RR			
# of YES	0	1	2
Probability	0.04	0.16 + 0.16	0.6

$\text{Var}[\# \text{ of YES}] = E[(X - \mu)^2] \approx 0.32$

- Variance: ( $\# \text{ people} = 2, p = 0.8$ )

$$\text{Var}[\hat{n}_{\text{YES}}] = \frac{\text{Var}[\# \text{ of YES}]}{(0.8 - 0.2)^2}$$

JRR		Better utility	
# of YES	0	1	2
Probability	0	0.2 + 0.2	0.6

$\text{Var}[\# \text{ of YES}] = E[(X - \mu)^2] = 0.24$

# Privacy: NOT as Simple as RR

- If any person can be an adversary

JRR: Joint distribution

	$T_1 = 1$	$T_1 = 0$
$T_2 = 1$	0.6	0.2
$T_2 = 0$	0.2	0

$T_1$ : I am an adversary 🤖

When I report untruthfully ( $T_1 = 0$ ),  
My partner will report truthfully ( $T_2 = 1$ )

# Privacy: NOT as Simple as RR

- If any person can be an adversary

JRR: Joint distribution

	$T_1 = 1$	$T_1 = 0$
$T_2 = 1$	0.6	0.2
$T_2 = 0$	0.2	0

$T_1$ : I am an adversary 🤖

When I report untruthfully ( $T_1 = 0$ ),  
My partner will report truthfully ( $T_2 = 1$ )

- Correlation results in privacy leakage (2 slides later)

# General JRR

- Correlated randomization with 2 persons  $x_{2i-1}$  and  $x_{2i}$

JRR: Joint distribution

	$T_{2i-1} = 1$	$T_{2i-1} = 0$
$T_{2i} = 1$	$p^2 + \rho pq$	$(1 - \rho)pq$
$T_{2i} = 0$	$(1 - \rho)pq$	$q^2 + \rho pq$



# General JRR

- Correlated randomization with 2 persons  $x_{2i-1}$  and  $x_{2i}$

JRR: Joint distribution

	$T_{2i-1} = 1$	$T_{2i-1} = 0$
$T_{2i} = 1$	$p^2 + \rho pq$	$(1 - \rho)pq$
$T_{2i} = 0$	$(1 - \rho)pq$	$q^2 + \rho pq$

$\rho \in [-1,1]$ :  
correlated coefficient



- RR is a special case of JRR with  $\rho = 0$  (no correlation)

# JRR – Privacy Model in This Paper

---

- Use random grouping to form 2-person groups for correlated randomization
  - secure shuffling by multi-party computing (MPC)

# JRR – Privacy Model in This Paper

---

- Use random grouping to form 2-person groups for correlated randomization
  - secure shuffling by multi-party computing (MPC)
- **Threat model:**
  - any person can be an adversary
  - if a group contains an adversary, the adversary knows **who is their partner** (after random grouping)

# JRR – Privacy Model in This Paper

- Use random grouping to form 2-person groups for correlated randomization
  - secure shuffling by multi-party computing (MPC)
- **Threat model:**
  - any person can be an adversary
  - if a group contains an adversary, the adversary knows **who is their partner** (after random grouping)

$$P \left[ \begin{array}{c} \text{the adversary} \\ \text{knows id} \end{array} \right] = \frac{m}{n-1}$$

→ *m*: # of adversaries

# JRR – Privacy Model in This Paper

- Use random grouping to form 2-person groups for correlated randomization
  - secure shuffling by multi-party computing (MPC)
- Threat model:**
  - any person can be an adversary
  - if a group contains an adversary, the adversary knows **who is their partner** (after random grouping)
  - the adversary cannot control randomness, but can **infer their partner's**

$$P[\text{the adversary knows id}] = \frac{m}{n-1}$$

$m$ : # of adversaries

$$P[\text{JRR}(x_2) = 1 \mid \text{JRR}(x_1) = 0] = \frac{(1-\rho)pq}{q} = (1-\rho)p$$

Adversary knows

Higher confidence

# JRR – Privacy Model in This Paper

- Use random grouping to form 2-person groups for correlated randomization
  - secure shuffling by multi-party computing (MPC)
- Threat model:**
  - any person can be an adversary
  - if a group contains an adversary, the adversary knows **who is their partner** (after random grouping)
  - the adversary cannot control randomness, but can **infer their partner's**

$$P[\text{the adversary knows id}] = \frac{m}{n-1} \quad \text{--- } m: \# \text{ of adversaries}$$

$$P[\text{JRR}(x_2) = 1 \mid \text{JRR}(x_1) = 0] = \frac{(1-\rho)pq}{q} = (1-\rho)p$$

Adversary knows

Higher confidence

Adversary's confidence

$$\frac{m}{n-1} \cdot (1-\rho)p$$

# JRR – Privacy Model in This Paper

- Use random grouping to form 2-person groups for correlated randomization

- secure shuffling by multi-party computing (MPC)

- Threat model:**

- any person can be an adversary
- if a group contains an adversary, the adversary can learn the result of the group
- the adversary cannot control randomization

$$p[\text{the adversary}] = \frac{m}{n-1} \rightarrow m: \# \text{ of adversaries}$$

**Privacy affected by**

$m$	# adversaries
$n$	# of persons
$\rho$	Correlated coefficient

(after random grouping)

$$P[\text{JRR}(x_2) = 1 \mid \text{JRR}(x_1) = 0] = \frac{(1-\rho)p}{q}$$

Adversary knows

Higher confidence

Adversary's confidence

$$\frac{m}{n-1} \cdot (1-\rho)p$$

# JRR – Formal Privacy & Utility

**Theorem.** Assume there is a set of data contributors  $\mathcal{T}_m$  whose reporting truthfulness is known to the adversary. For any data contributor  $i$ , the JRR mechanism satisfies:

$$\frac{\Pr[\text{JRR}(x_i) | \mathcal{T}_m]}{\Pr[\text{JRR}(x'_i) | \mathcal{T}_m]} \leq e^\varepsilon, \quad \text{where} \quad \varepsilon = \ln \frac{mp_{\max} + (n - m - 1)p}{mp_{\min} + (n - m - 1)q}.$$



# JRR – Formal Privacy & Utility

**Theorem.** Assume there is a set of data contributors  $\mathcal{T}_m$  whose reporting truthfulness is known to the adversary. For any data contributor  $i$ , the JRR mechanism satisfies:

$$\frac{\Pr[\text{JRR}(x_i) | \mathcal{T}_m]}{\Pr[\text{JRR}(x'_i) | \mathcal{T}_m]} \leq e^\varepsilon, \text{ where } \varepsilon = \ln \frac{mp_{\max} + (n - m - 1)p}{mp_{\min} + (n - m - 1)q}.$$

$m = |\mathcal{T}_m|$ :  
# of adversaries

$p_{\max} = \max\{(1 - \rho)p, p + \rho q\}$ :  
confidence of adversaries  
inferring a specific value

# JRR – Formal Privacy & Utility

**Theorem.** Assume there is a set of data contributors  $\mathcal{T}_m$  whose reporting truthfulness is known to the adversary. For any data contributor  $i$ , the JRR mechanism satisfies:

$$\frac{\Pr[\text{JRR}(x_i) | \mathcal{T}_m]}{\Pr[\text{JRR}(x'_i) | \mathcal{T}_m]} \leq e^\varepsilon, \text{ where } \varepsilon = \ln \frac{mp_{\max} + (n - m - 1)p}{mp_{\min} + (n - m - 1)q}.$$

**Theorem.** The variance of JRR's estimator  $\hat{n}_v$  is

$$\text{Var}[\hat{n}_v] = \frac{pq}{(p - q)^2} \cdot \left( n + \frac{\rho((2n_{\text{YES}} - n)^2 - n)}{n - 1} \right).$$

# JRR – Formal Privacy & Utility

**Theorem.** Assume there is a set of data contributors  $\mathcal{T}_m$  whose reporting truthfulness is known to the adversary. For any data contributor  $i$ , the JRR mechanism satisfies:

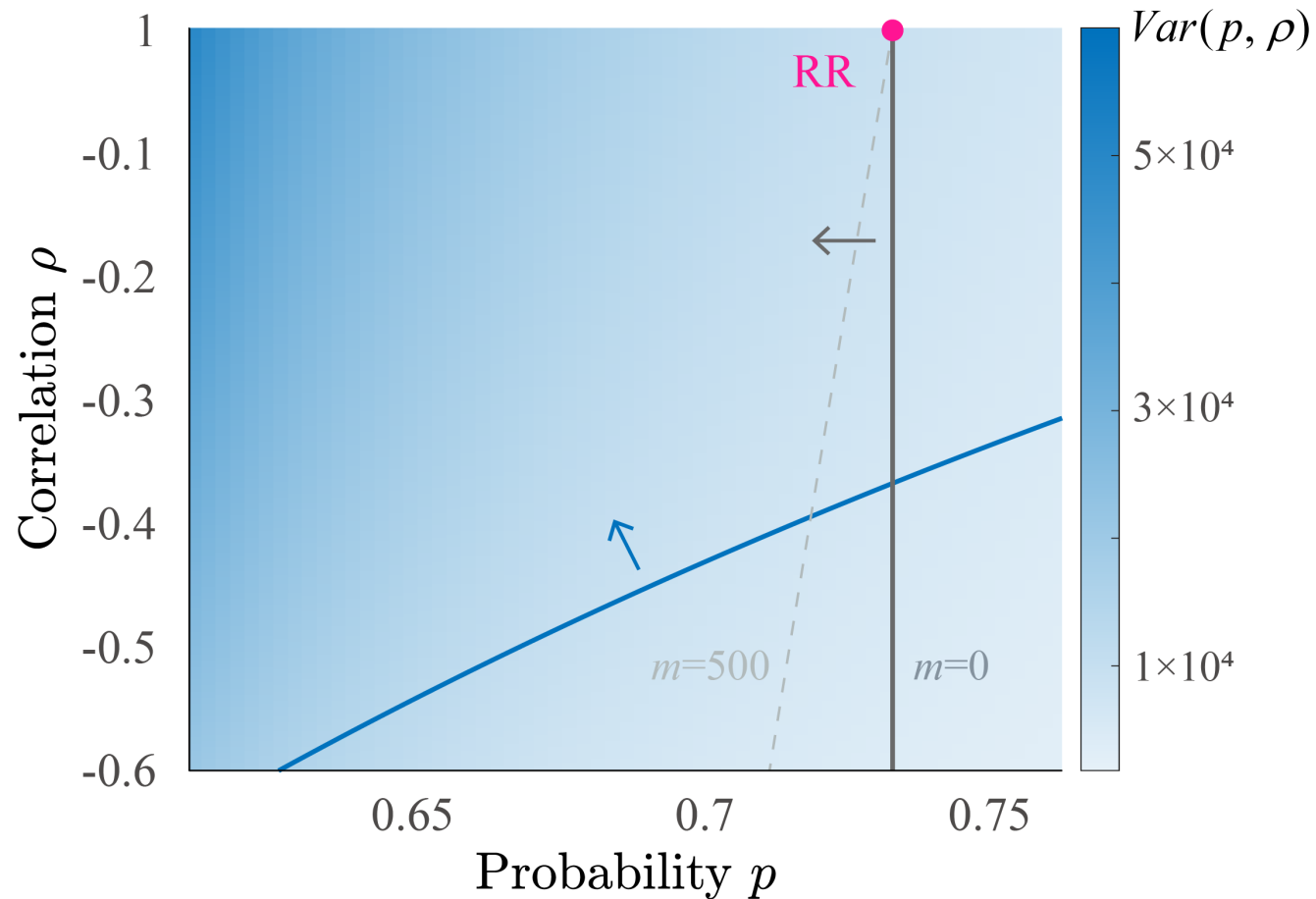
$$\frac{\Pr[\text{JRR}(x_i) | \mathcal{T}_m]}{\Pr[\text{JRR}(x'_i) | \mathcal{T}_m]} \leq e^\varepsilon, \text{ where } \varepsilon = \ln \frac{mp_{\max} + (n - m - 1)p}{mp_{\min} + (n - m - 1)q}.$$

**Theorem.** The variance of JRR's estimator  $\hat{n}_v$  is

$$\text{Var}[\hat{n}_v] = \frac{pq}{(p - q)^2} \cdot \left( \overset{\text{Independent}}{\underset{\text{randomization (RR)}}{n}} + \overset{\text{Affected by \# of original values}}{\overset{\text{Correlated}}{\underset{\text{randomization}}{\rho \left( \frac{(2n_{\text{YES}} - n)^2 - n}{n - 1} \right)}}} \right).$$

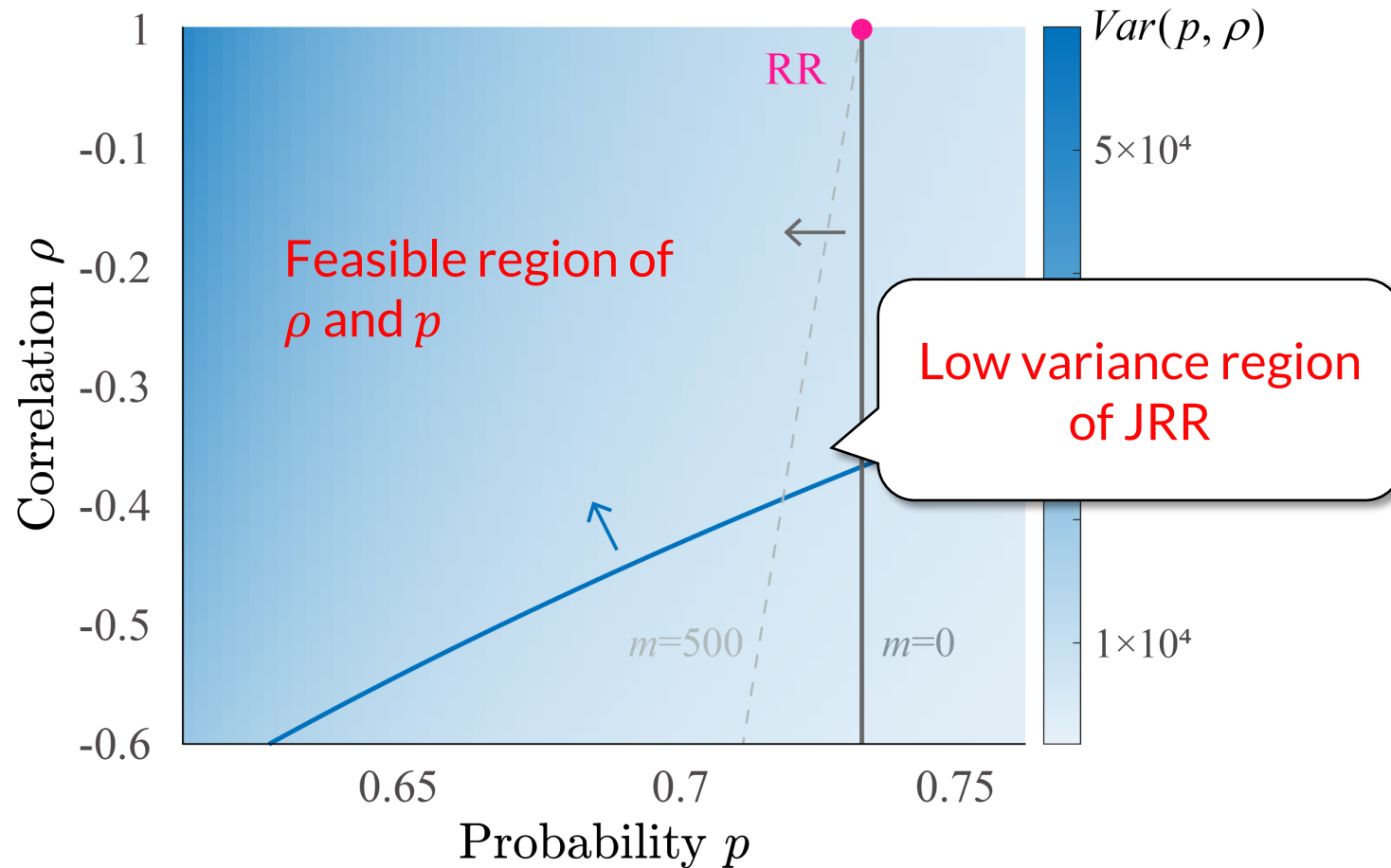
# JRR – Variance Heatmap

- Effect of  $\rho$ ,  $p$ , and  $m$  on variance (when  $\varepsilon = 1$ ,  $n = 10^4$ , and  $n_{\text{Yes}} = 200$ )



# JRR – Variance Heatmap

- Effect of  $\rho$ ,  $p$ , and  $m$  on variance (when  $\varepsilon = 1$ ,  $n = 10^4$ , and  $n_{\text{Yes}} = 200$ )



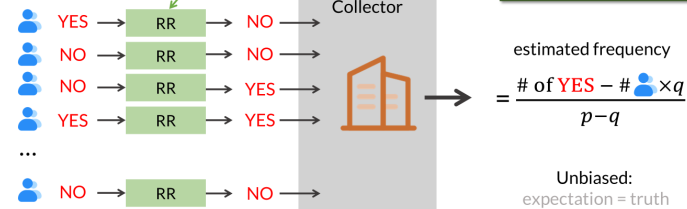
# Summary

- Correlated randomization can improve the data utility of frequency estimation
  - **JRR: Privacy & utility model for correlated randomization**
- 
- What's more in the paper
    - selection the tradeoff between  $\rho$  and  $p$
    - practical protocol design
    - prototype extensions to **non-binary data** and **larger-size group**
    - evaluations on synthetic and real-world datasets

# Locally Differentially Private Frequency Estimation via Joint Randomized Response

## Randomized Response for Privacy

- People **tend to lie** for such sensitive/embarrassing question - i.e. don't want to let the collector know
- Randomized Response: Randomize the truth **before answering**



RR: answer truth with probability  $p$

$$RR(x) = \begin{cases} x & \text{w.p. } p \\ \neg x & \text{w.p. } 1 - p \end{cases}$$

Ye Zheng

Locally Differentially Private Frequency Estimation via Joint Randomized Response

6

## This Paper: Joint RR (JRR)

- Joint randomization can boost data utility
- Example: 2-person ( $x_1 = \text{YES}$  and  $x_2 = \text{YES}$ ) with  $p = 0.8$  ( $P[T_1 = 1] = 0.8$ )

RR: Joint distribution

	$T_1 = 1$	$T_1 = 0$
$T_2 = 1$	0.64 ( $=0.8 \times 0.8$ )	0.16 ( $=0.2 \times 0.8$ )
$T_2 = 0$	0.16 ( $=0.8 \times 0.2$ )	0.04 ( $=0.2 \times 0.2$ )

Truthfulness of  $x_2$

Independent  $T_1$  and  $T_2$  ( $P[T_1 \cap T_2] = P[T_1] \cdot P[T_2]$ )

JRR: Joint distribution

	$T_1 = 1$	$T_1 = 0$
$T_2 = 1$	0.6	0.2
$T_2 = 0$	0.2	0

$P[T_1 = 0 \cap T_2 = 0] = 0$   
 $P[T_1 = 0] \cdot P[T_2 = 0] = 0.04$

NOT independent  $T_1$  and  $T_2$

$$P[T_1 = 1] = 0.6 + 0.2 = 0.8$$

Ye Zheng

Locally Differentially Private Frequency Estimation via Joint Randomized Response

17

## Utility: JRR's Variance

- Same estimator as RR

$$\text{Expectation: } E[\# \text{ of YES}] = \sum_{i=1}^n P_i$$

RR			
# of YES	0	1	2
Probability	0.04	0.16 + 0.16	0.6

$\text{Var}[\# \text{ of YES}] = E[(X - \mu)^2] \approx 0.32$

- Variance: ( $\# \text{ of people} = 2, p = 0.8$ )

$$\text{Var}[\hat{n}_{\text{YES}}] = \frac{\text{Var}[\# \text{ of YES}]}{(0.8 - 0.2)^2}$$

JRR Better utility			
# of YES	0	1	2
Probability	0	0.2 + 0.2	0.6

$\text{Var}[\# \text{ of YES}] = E[(X - \mu)^2] = 0.24$

Ye Zheng

Locally Differentially Private Frequency Estimation via Joint Randomized Response

21

## General JRR

- Correlated randomization with 2 persons  $x_{2i-1}$  and  $x_{2i}$

JRR: Joint distribution

	$T_{2i-1} = 1$	$T_{2i-1} = 0$
$T_{2i} = 1$	$p^2 + \rho pq$	$(1 - \rho)pq$
$T_{2i} = 0$	$(1 - \rho)pq$	$q^2 + \rho pq$

$\rho \in [-1, 1]$ : correlated coefficient

- RR is a special case of JRR with  $\rho = 0$  (no correlation)

Ye Zheng

Locally Differentially Private Frequency Estimation via Joint Randomized Response

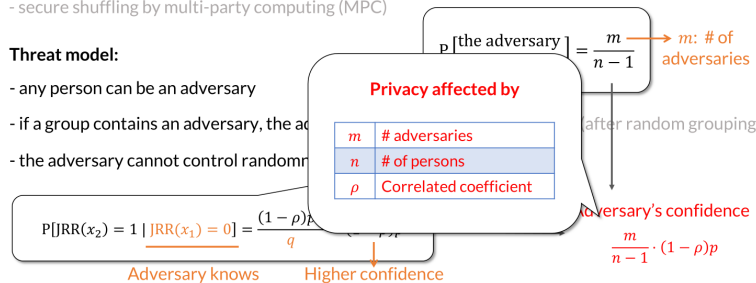
25

## JRR - Privacy Model in This Paper

- Use random grouping to form 2-person groups for correlated randomization - secure shuffling by multi-party computing (MPC)

- Threat model:

- any person can be an adversary
- if a group contains an adversary, the adversary can control the randomization
- the adversary cannot control randomization



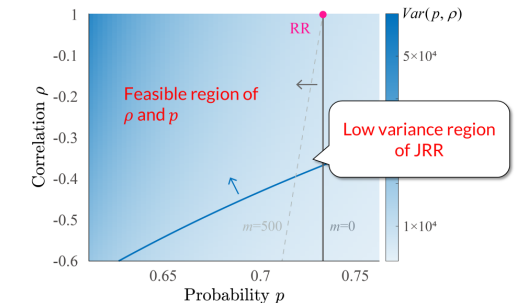
Ye Zheng

Locally Differentially Private Frequency Estimation via Joint Randomized Response

31

## JRR - Variance Heatmap

- Effect of  $\rho, p$ , and  $m$  on variance (when  $\epsilon = 1, n = 10^4$ , and  $n_{\text{YES}} = 200$ )



Ye Zheng

Locally Differentially Private Frequency Estimation via Joint Randomized Response

37

# Thank you!

