# (Local) Differential Privacy has NO Disparate Impact on Fairness

Héber H. Arcolezi, Karima Makhlouf, and Catuscia Palamidessi

Inria and École Polytechnique (IPP), Palaiseau, France
{heber.hwang-arcolezi,karima.makhlouf,catuscia}@lix.polytechnique.fr

DBSec, July 19th, 2023

# Motivation

# Differential Privacy (DP) and Fairness: Friends or Foes?

## Fairness Through Awareness

Cynthia Dwork
Microsoft Research S.V.
Mountain View, CA, USA
dwork@microsoft.com

Moritz Hardt[*]
IBM Research Almaden
San Jose, CA, USA
mhardt@us.ibm.com

Toniann Pitassi[†]
University of Toronto
Dept. of Computer Science
Toronto, ON, CANADA
toni@cs.toronto.edu

Omer Reingold
Microsoft Research S. V.
Mountain View, CA, USA
omer.reingold@microsoft.com

Richard Zemel[‡]
University of Toronto
Dept. of Computer Science
Toronto, ON, CANADA
zemel@cs.toronto.edu

**An Empirical Analysis of Fairness Notions under Differential Privacy[*]**

Anderson Santana de Oliveira,[1] Caelin Kaplan,[2] Khawla Mallat[1] Tanmay Chakraborty[3]

[1] SAP
[2] SAP and INRIA
[3] SAP and Eurecom
firstname.lastname@sap.com

## ON THE APPLICATION AND IMPACT OF $\epsilon$-DP AND FAIRNESS IN AMBULANCE ENGAGEMENT TIME PREDICTION

Selene Cerna & Catuscia Palamidessi
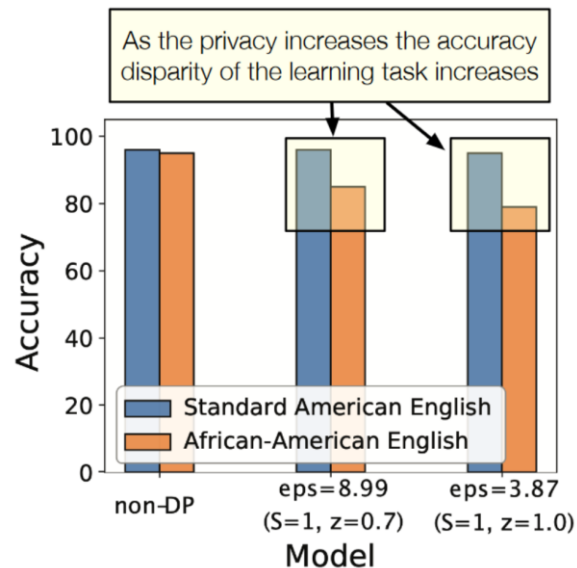Inria and École Polytechnique (IPP), Palaiseau, France
{selene-leya.cerna-nahuis, catuscia.palamidessi}@inria.fr

## Differential Privacy Has Disparate Impact on Model Accuracy

Eugene Bagdasaryan
Cornell Tech
eugene@cs.cornell.edu

Omid Poursaeed[*]
Cornell Tech
op63@cornell.edu

Vitaly Shmatikov
Cornell Tech
shmat@cs.cornell.edu

As the privacy increases the accuracy disparity of the learning task increases

# Differential Privacy (DP) and Fairness: Friends or Foes?

## Fairness Through Awareness

Cynthia Dwork
Microsoft Research S.V.
Mountain View, CA, USA
dwork@microsoft.com

Moritz Hardt[*]
IBM Research Almaden
San Jose, CA, USA
mhardt@us.ibm.com

Toniann Pitassi[†]
University of Toronto
Dept. of Computer Science
Toronto, ON, CANADA
toni@cs.toronto.edu

Omer Reingold
Microsoft Research S. V.
Mountain View, CA, USA
omer.reingold@microsoft.com

Richard Zemel[‡]
University of Toronto
Dept. of Computer Science
Toronto, ON, CANADA
zemel@cs.toronto.edu

## An Empirical Analysis of Fairness Notions under Differential Privacy[*]

Anderson Santana de Oliveira,[1] Caelin Kaplan,[2] Khawla Mallat[1] Tanmay Chakraborty[3]

[1] SAP
[2] SAP and INRIA
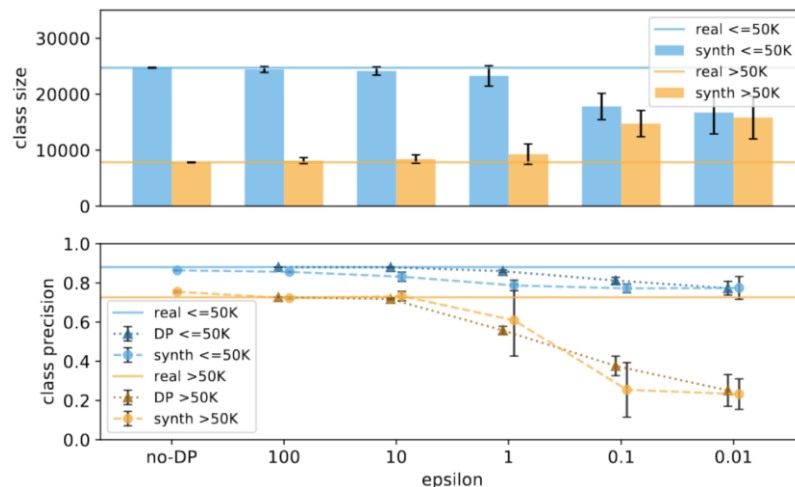[3] SAP and Eurecom
firstname.lastname@sap.com

## ON THE APPLICATION AND IMPACT OF $\epsilon$-DP AND FAIRNESS IN AMBULANCE ENGAGEMENT TIME PREDICTION

Selene Cerna & Catuscia Palamidessi
Inria and École Polytechnique (IPP), Palaiseau, France
{selene-leya.cerna-nahuis, catuscia.palamidessi}@inria.fr

## Robin Hood and Matthew Effects: Differential Privacy Has Disparate Impact on Synthetic Data
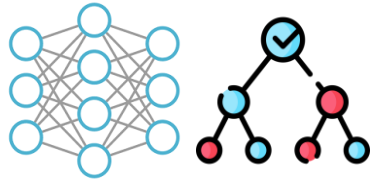
Georgi Ganev[1,2]   Bristena Oprisanu[1]   Emiliano De Cristofaro[1]
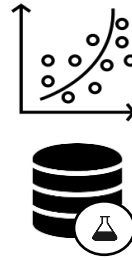
# Local DP (LDP) and Fairness: Friends or Foes?
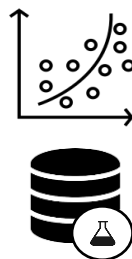


Dataset after LDP      ML algorithm      Output      Fairness metrics impacted?
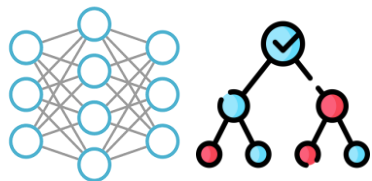
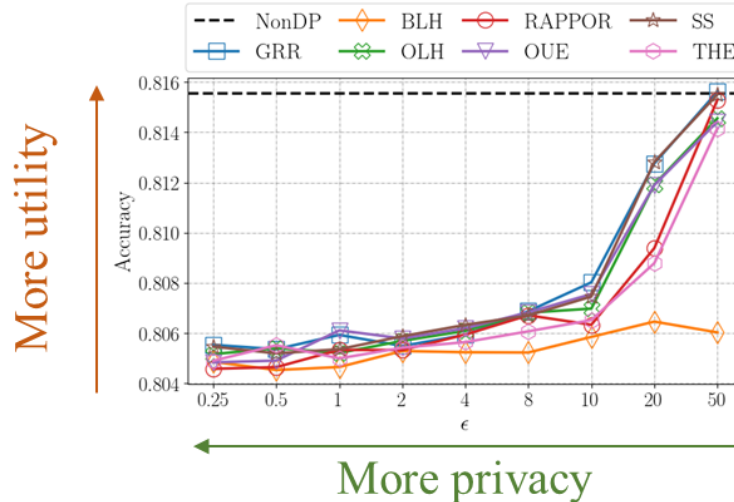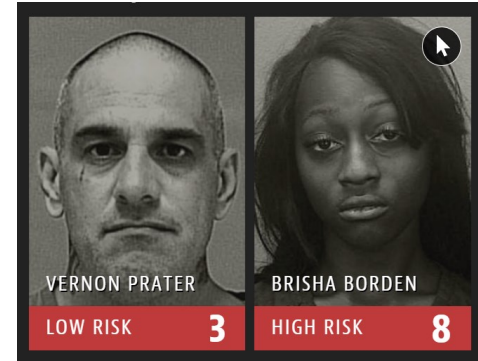# Local DP (LDP) and Fairness: Friends! ~~or~~ ~~Foes~~?



Dataset after LDP          ML algorithm          Output          Fairness metrics impacted?



More Fairness — More privacy



More utility — More privacy

# Differential Privacy (DP) and Fairness: Friends or Foes?

| Paper | Task | Privacy | Details | Results |
|---|---|---|---|---|
| DP Has Disparate Impact on Model Accuracy (NeurIPS 2019) | Classification | Central DP | DP-SGD w/ same hyperparameters as the non-private baseline. | Foes |
| Robin Hood and Matthew Effects: DP Has Disparate Impact on Synthetic Data (ICML 2022) | Synthetic data generation + classification | Central DP | DP generative models w/ same hyperparameters as the non-private baseline. | Foes |
| An Empirical Analysis of Fairness Notions under DP (PPAI 2023) | Classification | Central DP | DP-SGD: search for optimal hyperparameters. | Minor impact |
| DP has Bounded Impact on Fairness in Classification (ICML 2023) | Classification | Central DP | DP-SGD: Theory. | Bounded impact |
| Fair Learning with Private Demographic Data (ICML 2020) | Classification | Local DP | LDP on single attribute + fairness mitigation mechanism. | |
| On the application and impact of $\epsilon$-DP and fairness in ambulance engagement time prediction (ICLR 2023) | Classification | Local DP | LDP on multiple attributes. | Friends |
| **Our (DBSec 2023)** | Classification | Local DP | LDP on multiple attributes. | Friends |

Inria    ÉCOLE POLYTECHNIQUE
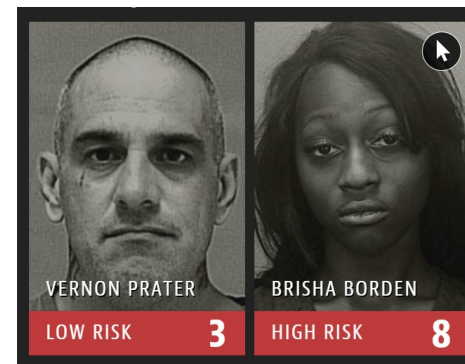
# Outline

# Fairness Metrics

**Fairness** [Cambridge Dictionary]: The quality of treating people equally or in a way that is right or reasonable.
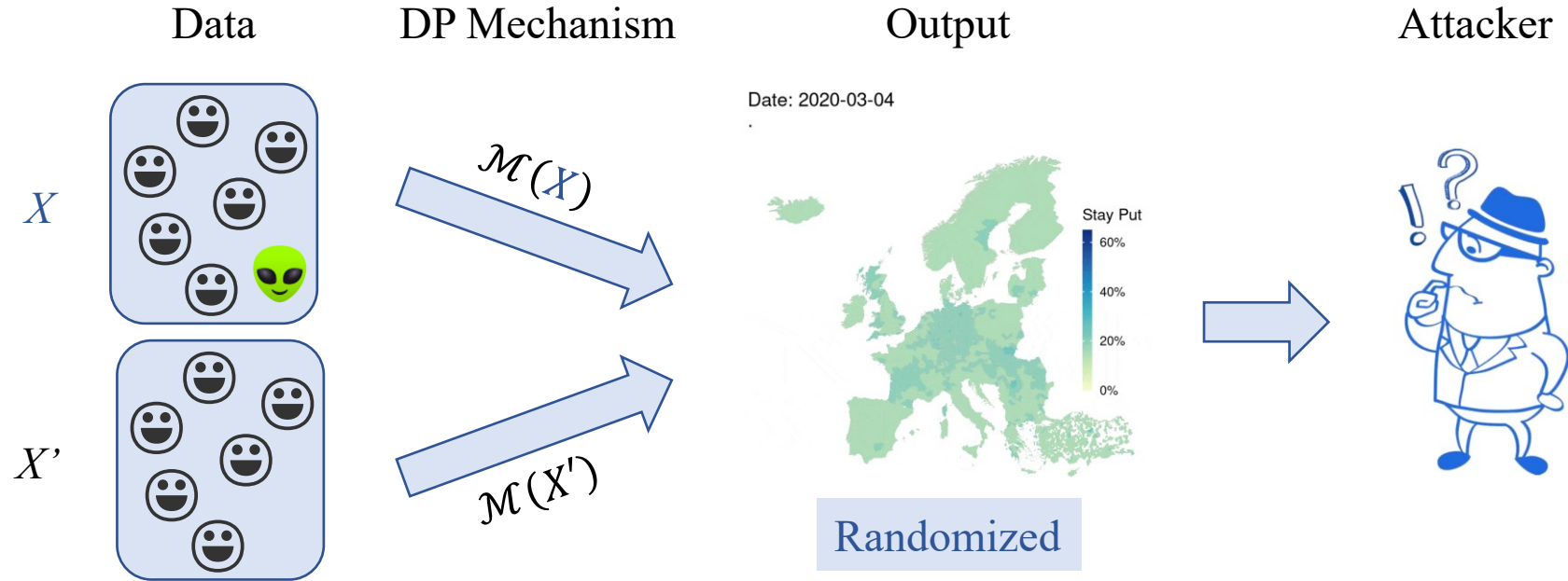
# Fairness Metrics

**Fairness** [Cambridge Dictionary]: The quality of treating people equally or in a way that is right or reasonable.

Protected attribute: $A_p \in \{0,1\}$
Target, predictor: $Y, \hat{Y} \in \{0,1\}$



VERNON PRATER  BRISHA BORDEN
LOW RISK  **3**  HIGH RISK  **8**

| Fairness Metric | Equation | When Satisfied? |
|---|---|---|
| Disparate Impact (DI) | $$\frac{\Pr[\hat{Y}=1|A_p=0]}{\Pr[\hat{Y}=1|A_p=1]}$$ | 1 |
| Statistical Parity Difference (SPD) | $\Pr[\hat{Y}=1|A_p=1] - \Pr[\hat{Y}=1|A_p=0]$ | 0 |
| Equal Opportunity Difference (EOD) | $\Pr[\hat{Y}=1|Y=1,A_p=1] - \Pr[\hat{Y}=1|Y=1,A_p=0]$ | 0 |
| Overall Accuracy Difference (OAD) | $\Pr[\hat{Y}=Y|A_p=1] - \Pr[\hat{Y}=Y|A_p=0]$ | 0 |

Inria  ÉCOLE POLYTECHNIQUE

# Differential Privacy (DP) [Dwork et al, 2006]



The attacker **cannot** tell if 👽 is in the sample

# Differential Privacy (DP) [Dwork et al, 2006; Duchi et al, 2013]



**Centralized DP:**

✅ High utility.

❌ Need to trust the server.

❌❌ **Data breaches, data misuse, etc.**

**Local DP (LDP):**

✅ No need to trust the server.

❌ Low utility.

# LDP: Formal Definition & Properties [Duchi et al, 2013]

*Def (ε-LDP)*. A randomized mechanism $\mathcal{M}$ satisfies $\epsilon$-LDP, where $\epsilon \geq 0$, if for **any two inputs** $v, v' \in \text{Domain}(\mathcal{M})$ and for **any output** $z \in \text{Range}(\mathcal{M})$:

$$\frac{\Pr[\mathcal{M}(v) = z]}{\Pr[\mathcal{M}(v') = z]} \leq e^{\epsilon}$$

Privacy Loss

Utility — Privacy

# LDP: Formal Definition & Properties [Duchi et al, 2013]

*Def ($\epsilon$-LDP).* A randomized mechanism $\mathcal{M}$ satisfies $\epsilon$-LDP, where $\epsilon \geq 0$, if for **any two inputs** $v, v' \in \text{Domain}(\mathcal{M})$ and for **any output** $z \in \text{Range}(\mathcal{M})$:

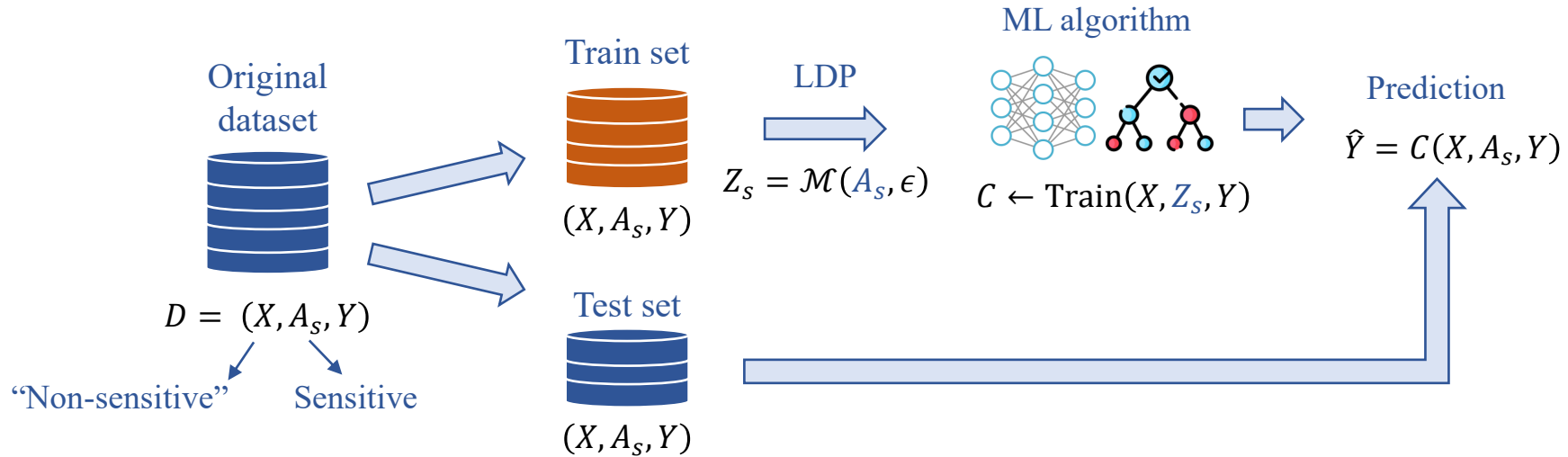$$\frac{\Pr[\mathcal{M}(v) = z]}{\Pr[\mathcal{M}(v') = z]} \leq e^{\epsilon}$$

Privacy Loss

Utility    Privacy

Fundamental (L)DP properties [Dwork et al, 2006]:

- **Post-processing** $\rightarrow$ if $\mathcal{M}$ is $\epsilon$-LDP, then the composition $f(\mathcal{M})$ is $\epsilon$-LDP for any $f$.

- **Composition** $\rightarrow$ Let $\mathcal{M}_1$ be a $\epsilon_1$-LDP mechanism and $\mathcal{M}_2$ a $\epsilon_2$-LDP mechanism. Then, the composed mechanism $\mathcal{M} = \big(\mathcal{M}_1(v), \mathcal{M}_2(v)\big)$ is $(\epsilon_1 + \epsilon_2)$-LDP.

# Outline

# Problem Statement



ML algorithm

Train set

Original
dataset

LDP

Prediction

$Z_s = \mathcal{M}(A_s, \epsilon)$

$C \leftarrow \text{Train}(X, Z_s, Y)$

$\hat{Y} = C(X, A_s, Y)$

$(X, A_s, Y)$

$D = (X, A_s, Y)$

Test set

"Non-sensitive"     Sensitive

$(X, A_s, Y)$

User's goal:

- Sanitize multiple sensitive attributes ($|A_s| \geq 2$) independently with $\epsilon$-LDP.

Server's goal:

- Train a Machine Learning (ML) classifier on sanitized data $(X, Z_s, Y)$.

# Research Questions (RQs) & Assumptions

- **RQ1:** How does LDP pre-processing impacts fairness & utility?

- **RQ2:** How to better split the privacy budget $\epsilon$ for $d_s = |A_s|$ sensitive attributes?

- **RQ3:** Which LDP protocol lead to the best privacy-utility-fairness trade-off?

# Research Questions (RQs) & Assumptions

- RQ1: How does LDP pre-processing impacts fairness & utility?

  - (Fairness) protected attribute $A_p$ is always a sensitive attribute $A_p \in A_s$;

  - Empirical results w/ 3 datasets, 4 fairness metrics, and 4 utility metrics.

- RQ2: How to better split the privacy budget $\epsilon$ for $d_s = |A_s|$ sensitive attributes?

- RQ3: Which LDP protocol lead to the best privacy-utility-fairness trade-off?

# Research Questions (RQs) & Assumptions

- RQ1: How does LDP pre-processing impacts fairness & utility?

  - (Fairness) protected attribute $A_p$ is always a sensitive attribute $A_p \in A_s$;

  - Empirical results w/ 3 datasets, 4 fairness metrics, and 4 utility metrics.

- RQ2: How to better split the privacy budget $\epsilon$ for $d_s = |A_s|$ sensitive attributes?

  - State-of-the-art: Uniform splitting $\rightarrow \epsilon_j = \frac{\epsilon}{d_s}$ for $j \in A_s$;

  - Our solution: $k$-based $\rightarrow \epsilon_j = \frac{\epsilon \cdot k_j}{\sum_{i=1}^{d_s} k_i}$ for $j \in A_s$, $k_j = |A_j|$.

  > $\epsilon$-LDP following the sequential composition

- RQ3: Which LDP protocol lead to the best privacy-utility-fairness trade-off?
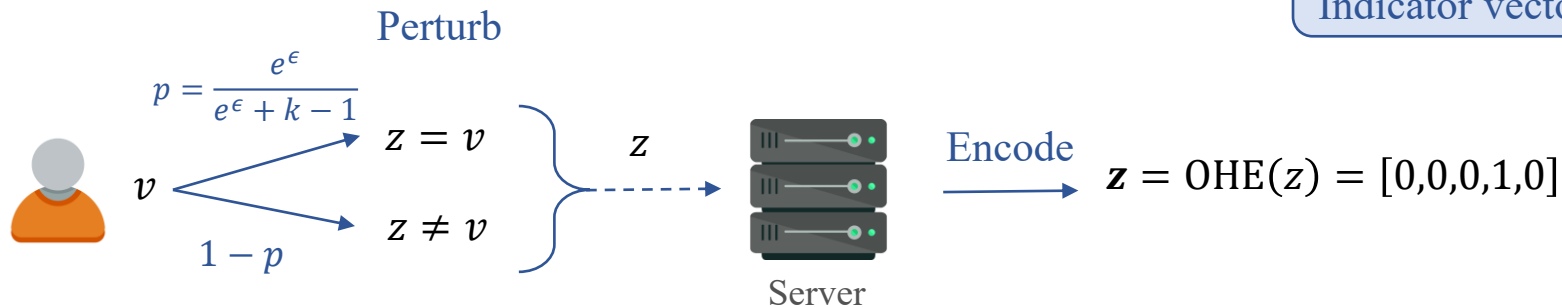
# Research Questions (RQs) & Assumptions

- RQ1: How does LDP pre-processing impacts fairness & utility?

  - (Fairness) protected attribute $A_p$ is always a sensitive attribute $A_p \in A_s$;

  - Empirical results w/ 3 datasets, 4 fairness metrics, and 4 utility metrics.

- RQ2: How to better split the privacy budget $\epsilon$ for $d_s = |A_s|$ sensitive attributes?

  - State-of-the-art: Uniform splitting $\rightarrow \epsilon_j = \frac{\epsilon}{d_s}$ for $j \in A_s$;

  - Our solution: $k$-based $\rightarrow \epsilon_j = \frac{\epsilon \cdot k_j}{\sum_{i=1}^{d_s} k_i}$ for $j \in A_s$, $k_j = |A_j|$.

  $\epsilon$-LDP following the sequential composition

- RQ3: Which LDP protocol lead to the best privacy-utility-fairness trade-off?

  - Benchmarked 7 state-of-the-art LDP protocols;

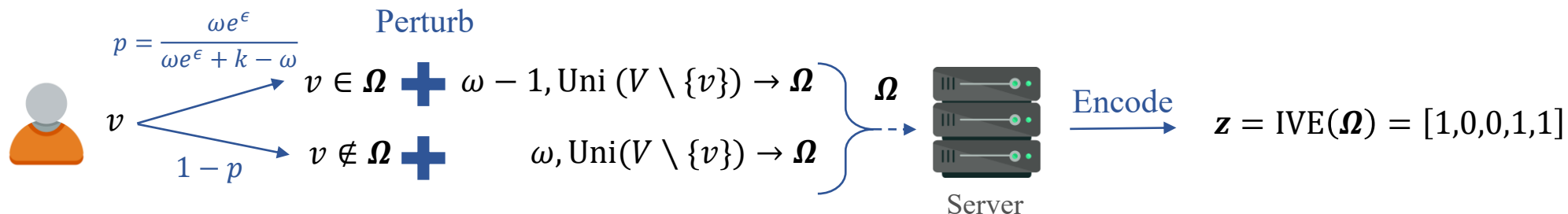  - Post-processed $\epsilon$-LDP report for "homogeneous encoding" at the server side.

# LDP Protocols & Server's "Homogeneous" Encoding

**Generalized Randomized Response (GRR)**
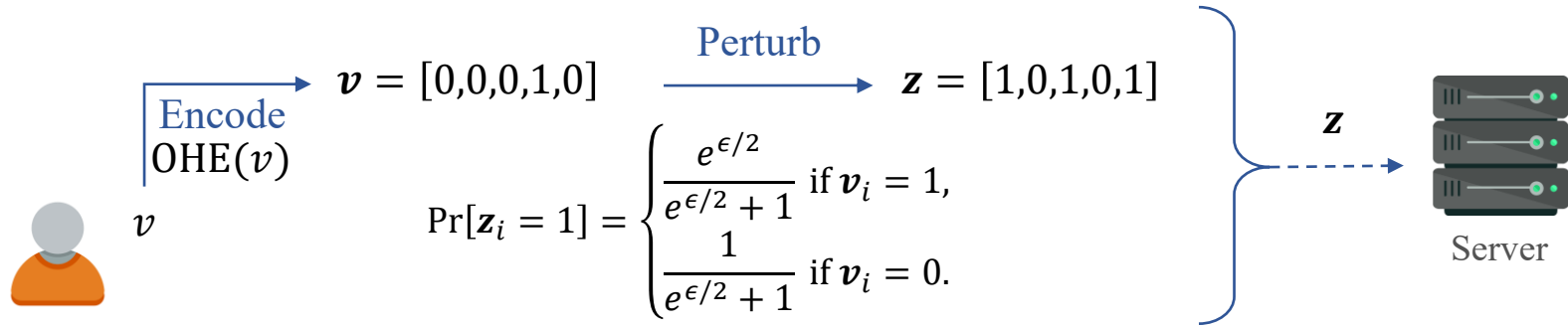
One-hot-encoding (OHE)
Indicator vector encoding (IVE)

Perturb

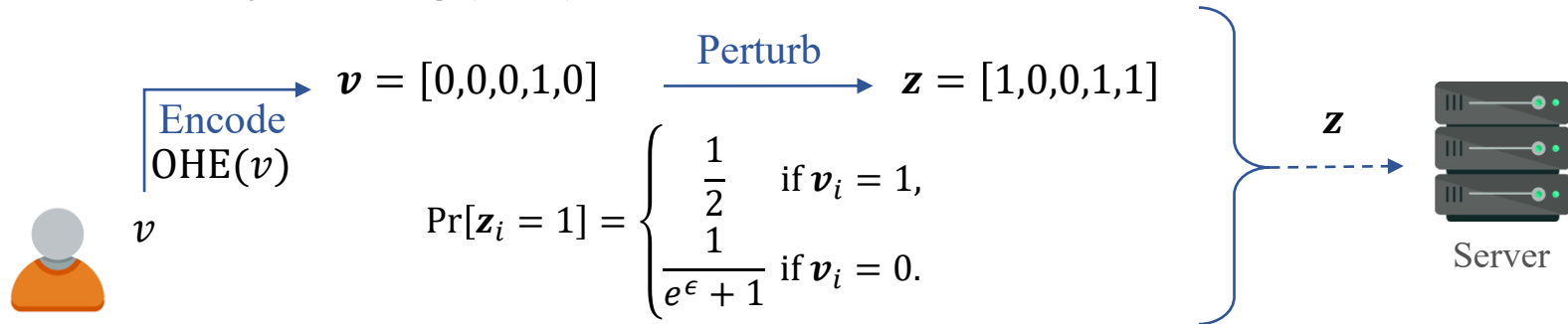$$p = \frac{e^{\epsilon}}{e^{\epsilon} + k - 1}$$

$v$ → $z = v$

$z \neq v$

$1 - p$

$z$ → Server

Encode → $\mathbf{z} = \text{OHE}(z) = [0,0,0,1,0]$

**Subset Selection (SS)**

Perturb

$$p = \frac{\omega e^{\epsilon}}{\omega e^{\epsilon} + k - \omega}$$

$v$ → $v \in \mathbf{\Omega}$ ✚ $\omega - 1, \text{Uni}(V \setminus \{v\}) \to \mathbf{\Omega}$

$v \notin \mathbf{\Omega}$ ✚ $\omega, \text{Uni}(V \setminus \{v\}) \to \mathbf{\Omega}$

$1 - p$

$\mathbf{\Omega}$ → Server

Encode → $\mathbf{z} = \text{IVE}(\mathbf{\Omega}) = [1,0,0,1,1]$

# LDP Protocols & Server's "Homogeneous" Encoding

**RAPPOR**

$$v = [0,0,0,1,0] \xrightarrow{\text{Perturb}} z = [1,0,1,0,1]$$

Encode
$\text{OHE}(v)$

$v$

$$\Pr[z_i = 1] = \begin{cases} \dfrac{e^{\epsilon/2}}{e^{\epsilon/2} + 1} & \text{if } v_i = 1, \\ \dfrac{1}{e^{\epsilon/2} + 1} & \text{if } v_i = 0. \end{cases}$$

$z$

Server

**Optimized Unary Encoding (OUE)**

$$v = [0,0,0,1,0] \xrightarrow{\text{Perturb}} z = [1,0,0,1,1]$$

Encode
$\text{OHE}(v)$

$v$

$$\Pr[z_i = 1] = \begin{cases} \dfrac{1}{2} & \text{if } v_i = 1, \\ \dfrac{1}{e^{\epsilon} + 1} & \text{if } v_i = 0. \end{cases}$$

$z$

Server

# LDP Protocols & Server's "Homogeneous" Encoding

**Local Hashing (LH)**

$$\boxed{\begin{array}{l}\text{Binary LH: } g = 2 \\ \text{Optimal LH: } g = e^\epsilon + 1\end{array}}$$



**Thresholding w/ Histogram Encoding (THE)**

# Outline

# Setting of Experiments

Three datasets:

- Adult, ACSCoverage, LSAC.

Four fairness metrics:

- DI, SPD, EOD, AOD.

ML Classifier:

- LGBM w/ fixed hyperparameters;

- Train/test split as 80/20.

Seven LDP protocols:

- GRR, SS, RAPPOR, OUE, BLH, OLH, THE.

Two privacy budget splitting solutions:

- Uniform and $k$-based.

Fixed $|A_s| = 4$

| Dataset | $n$ | $A_p$ | $A_s$, domain size $k$ | $Y$ |
|---------|-----|-------|------------------------|-----|
| Adult | 45849 | gender | - gender, $k = 2$<br>- race, $k = 5$<br>- native country, $k = 41$<br>- age, $k = 74$ | income |
| ACSCoverage | 98739 | DIS | - DIS, $k = 2$<br>- AGEP, $k = 50$<br>- SEX, $k = 2$<br>- SCHL, $k = 24$ | PUBCOV |
| LSAC | 20427 | race | - race, $k = 2$<br>- gender, $k = 2$<br>- family income, $k = 5$<br>- full time, $k = 2$ | pass bar |

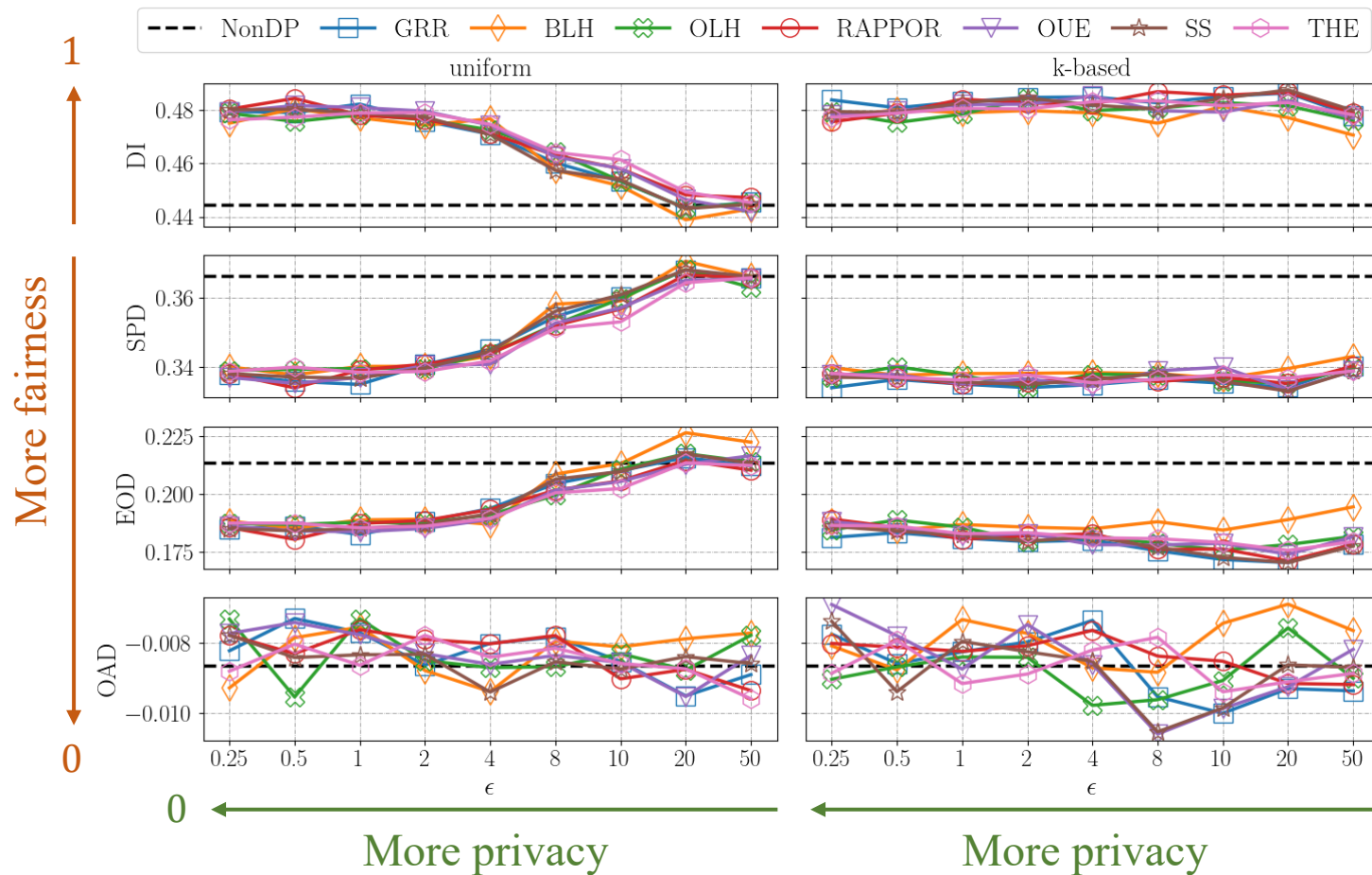Stability: average over 20 runs

Inria ÉCOLE POLYTECHNIQUE

# Impact of LDP on Fairness

$$DI = \frac{\Pr[\hat{Y} = 1 | A_p = 0]}{\Pr[\hat{Y} = 1 | A_p = 1]} \rightarrow 1$$

$$SPD = \Pr[\hat{Y} = 1 | A_p = 1]$$
$$-\Pr[\hat{Y} = 1 | A_p = 0] \rightarrow 0$$

$$EOD = \Pr[\hat{Y} = 1 | Y = 1, A_p = 1]$$
$$-\Pr[\hat{Y} = 1 | Y = 1, A_p = 0] \rightarrow 0$$

$$OAD = \Pr[\hat{Y} = Y | A_p = 1] -$$
$$\Pr[\hat{Y} = Y | A_p = 0] \rightarrow 0$$

# Impact of LDP on Fairness
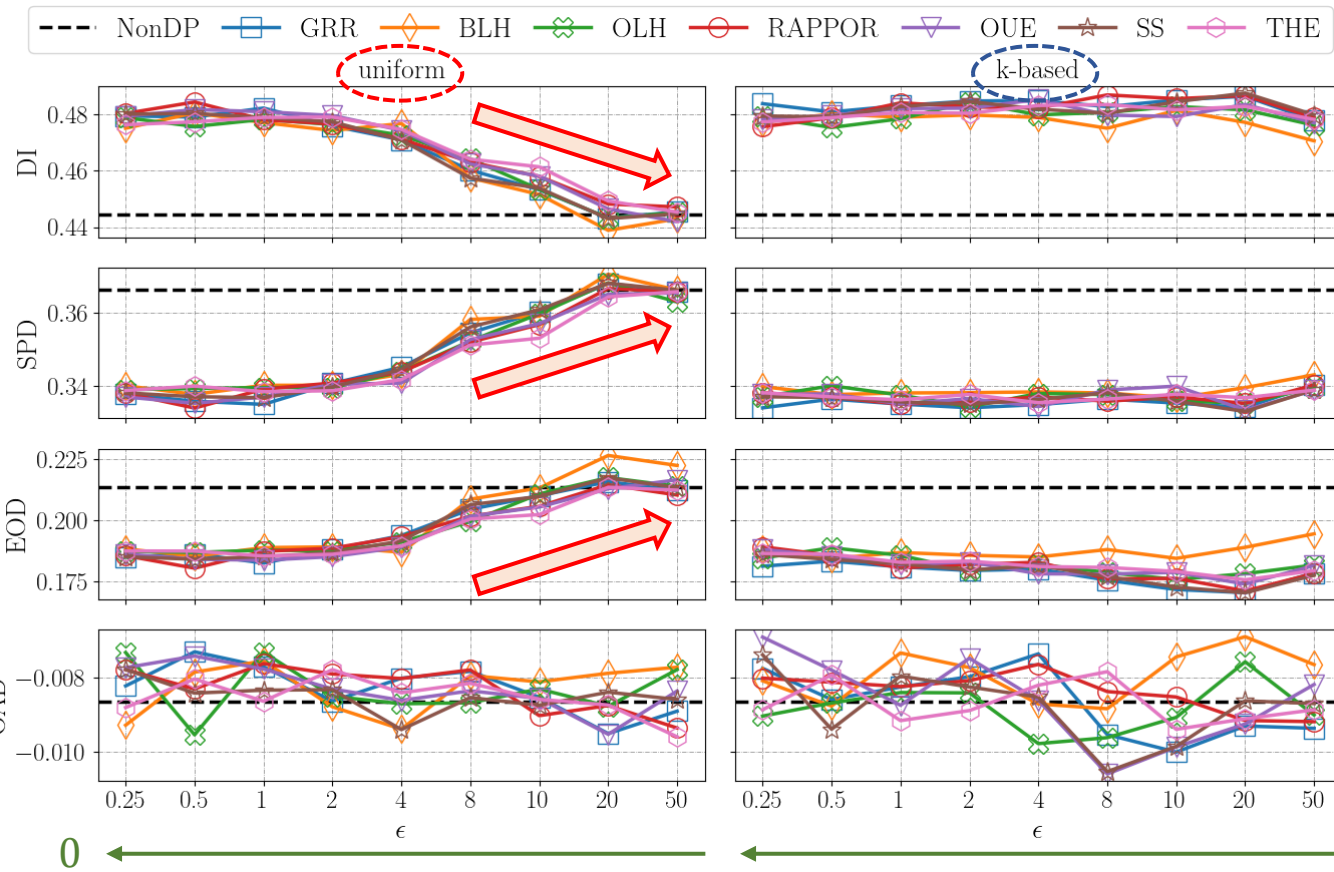
Uniform: goes towards the 'bad' baseline fairness metrics

$$DI = \frac{\Pr[\hat{Y} = 1 | A_p = 0]}{\Pr[\hat{Y} = 1 | A_p = 1]} \rightarrow 1$$

$$SPD = \Pr[\hat{Y} = 1 | A_p = 1] - \Pr[\hat{Y} = 1 | A_p = 0] \rightarrow 0$$

$$EOD = \Pr[\hat{Y} = 1 | Y = 1, A_p = 1] - \Pr[\hat{Y} = 1 | Y = 1, A_p = 0] \rightarrow 0$$

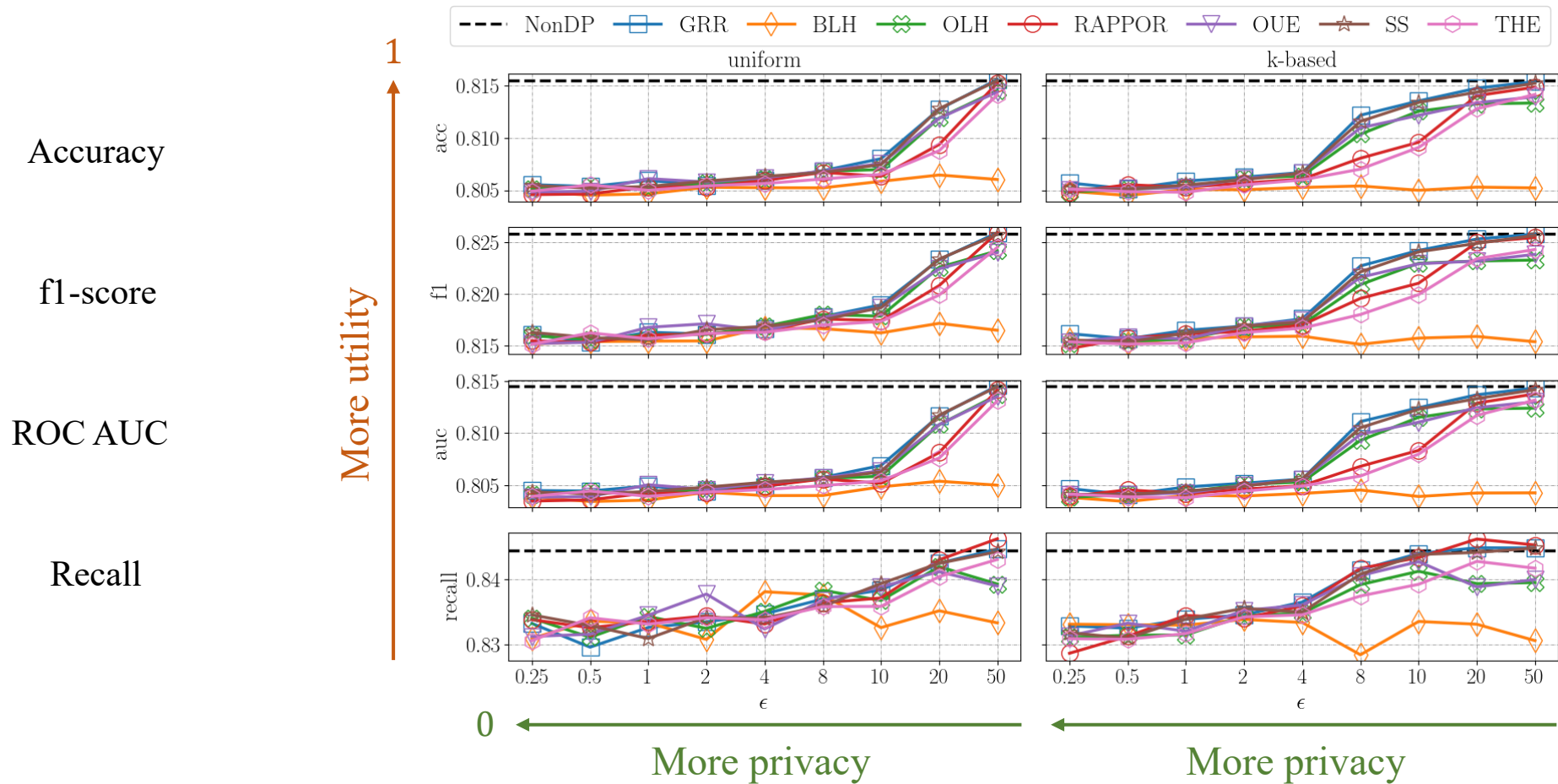$$OAD = \Pr[\hat{Y} = Y | A_p = 1] - \Pr[\hat{Y} = Y | A_p = 0] \rightarrow 0$$
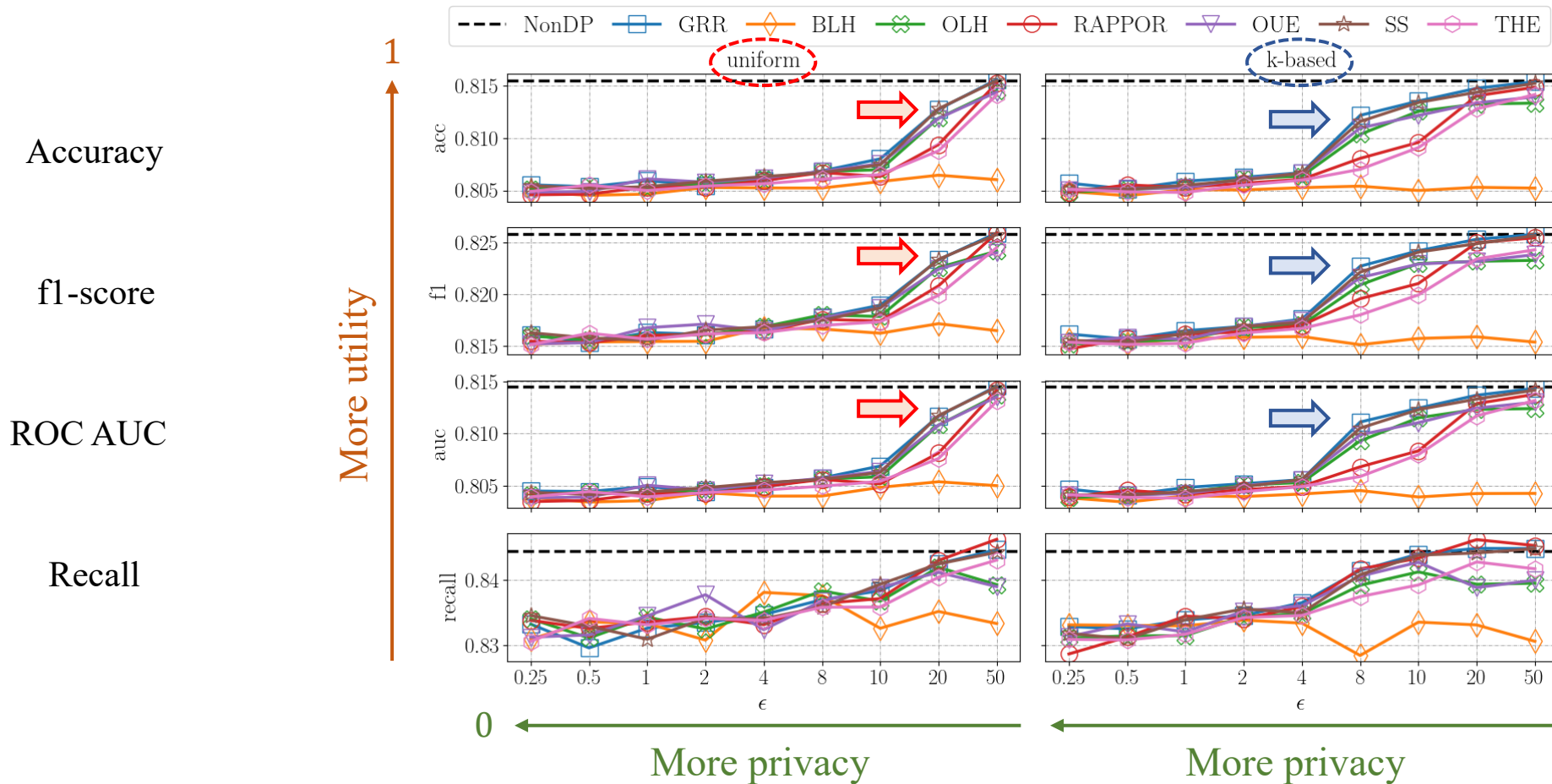
More fairness
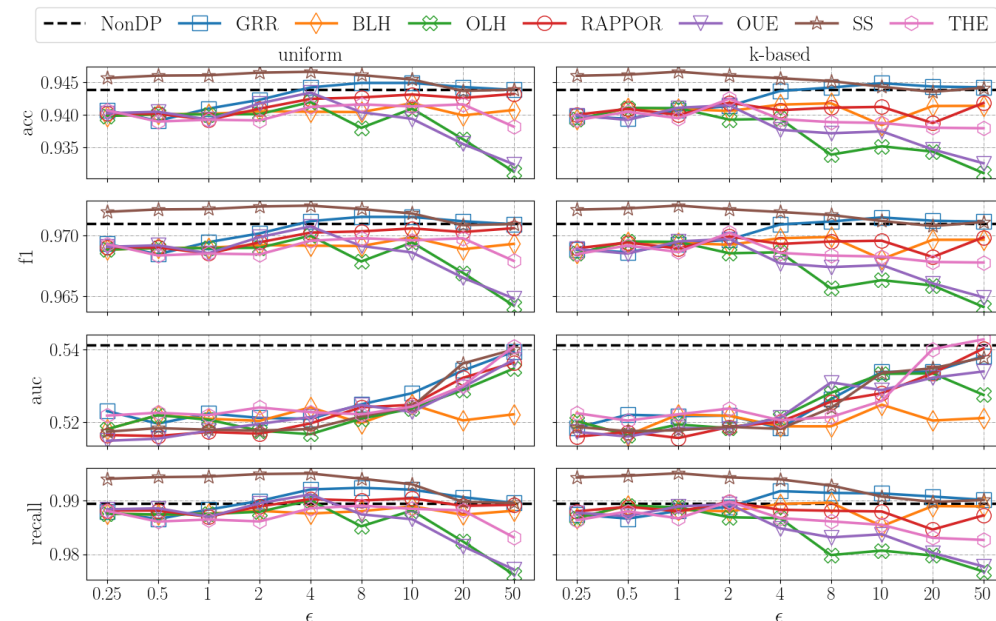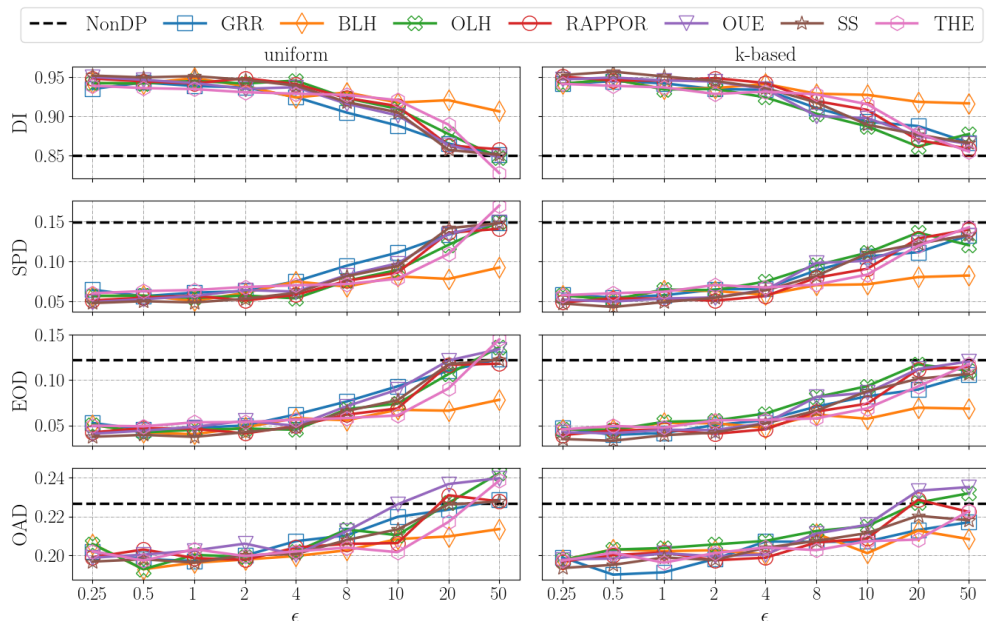
# Impact of LDP on Utility

# Impact of LDP on Utility



*k*-based: approaches faster the 'good' baseline utility metrics

# Impact of LDP on Fairness & Utility: Generic? → Yes!

Appendix Experiments: $|A_s| = \text{Uniform}([2, 6])$.



More privacy

More privacy

# Outline

# Takeaway Messages

**Conclusions:**

- DP does not necessarily lead to worsened fairness in ML;

- (L)DP pre-processing positively affects fairness w/ minor utility impact;

- Our $k$-based solution leads to better privacy-utility-fairness trade-off;

- Mechanism w/ best privacy-utility-fairness trade-off: GRR and SS.

# Takeaway Messages

**Conclusions:**

- DP does not necessarely lead to worsened fairness in ML;

- (L)DP pre-processing positively affects fairness w/ minor utility impact;

- Our $k$-based solution leads to better privacy-utility-fairness trade-off;

- Mechanism w/ best privacy-utility-fairness trade-off: GRR and SS.

**Perspectives:**

- Formalize our findings (*i.e.*, LDP & fairness trade-off);

- Introduce optimal mechanisms for privacy-fairness-aware ML;

- Study the impact of LDP pre-processing on different ML algorithms.

# (Local) Differential Privacy has NO Disparate Impact on Fairness

Héber H. Arcolezi, Karima Makhlouf, and Catuscia Palamidessi

Inria and École Polytechnique (IPP), Palaiseau, France

PAPER

ARTIFACT

CONTACT

hharcolezi.github.io          heber.hwang-arcolezi@inria.fr          @hharcolezi