

# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection

**Yuxin Wang, Duanyu Feng, Yongfu Dai, Zhengyu Chen, Jimin Huang,  
Sophia Ananiadou, Qianqian Xie, Hao Wang**



四川大學  
SICHUAN UNIVERSITY



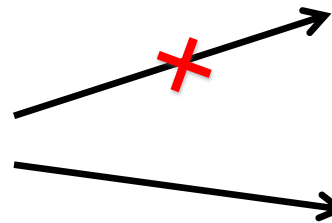
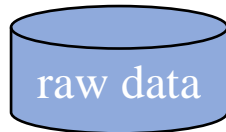
# Data Synthesis (What is/for)

- Synthetic data is defined as data obtained from a generative process that **learns the properties of the original data**.
- It intends to synthesize new samples that are related to but can **not be mapped back to the original data**.
- The generation of synthetic data can be used for **anonymization, regularization, oversampling, semi-supervised learning**, and several other tasks.

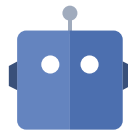
Sorry, we can't give it to you directly.



Data Owner



OK, I see...

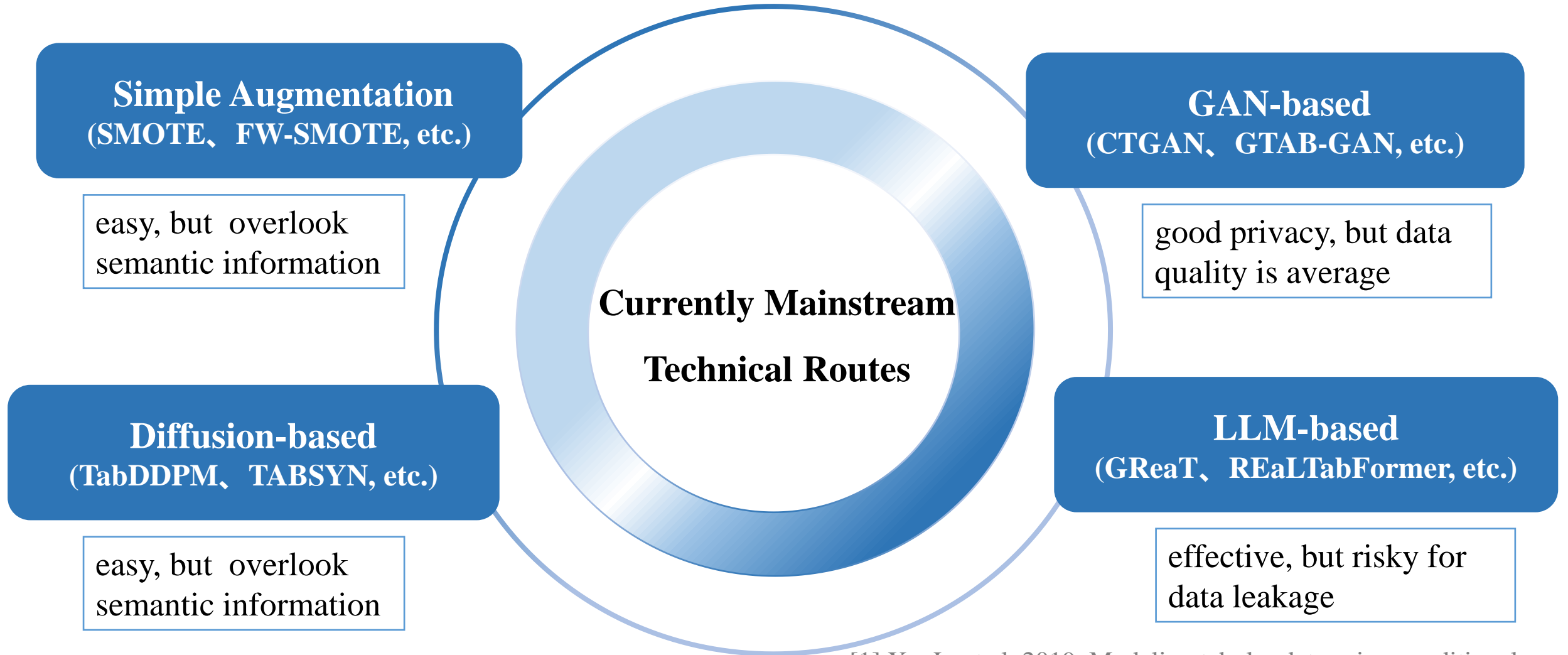


Synthetic data can be a good substitute.

[1] Assefa, S. A. et al. 2022. Generating synthetic data in finance: opportunities, challenges and pitfalls.

[2] Fonseca, J. et al. 2023. Tabular and latent space synthetic data generation: a literature review.

# Data Synthesis (How to do)



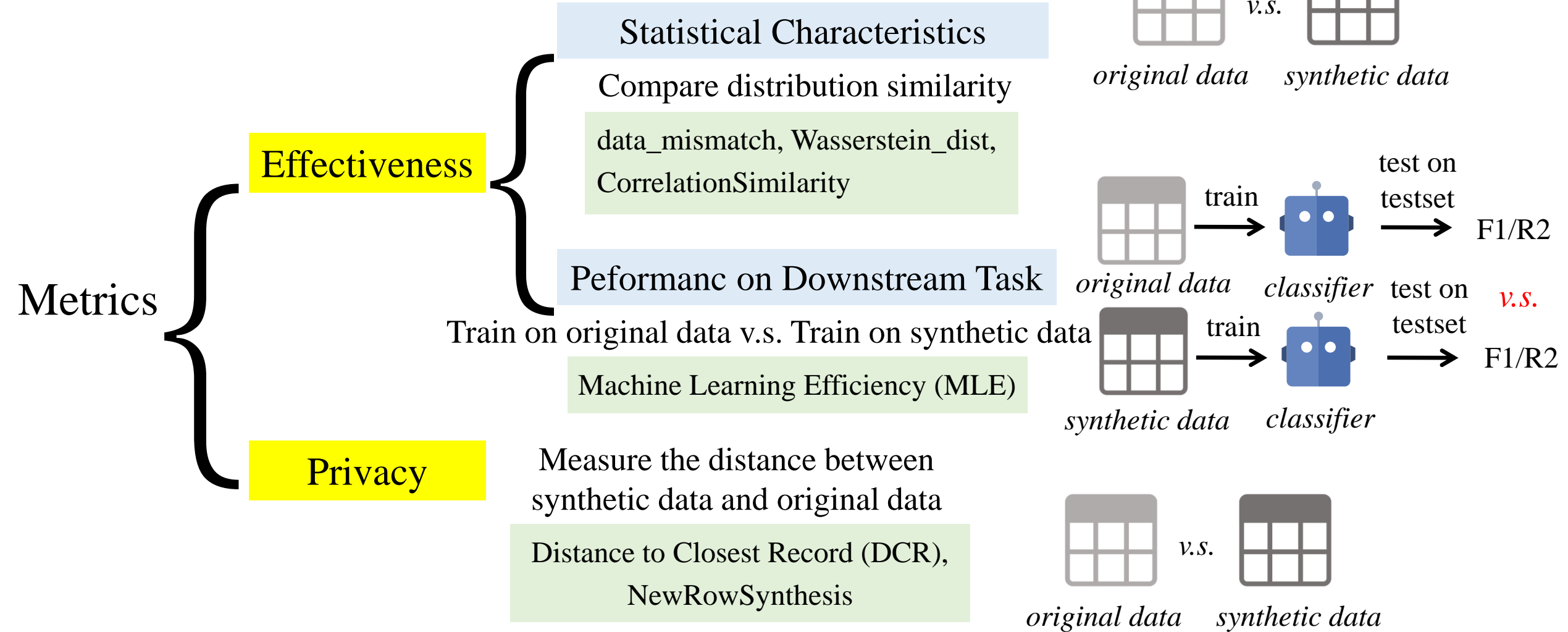
[1] Xu, L. et al. 2019. Modeling tabular data using conditional gan.

[2] Chawla, N. V. et al. 2002. SMOTE: synthetic minority over-sampling technique.

[3] Zhang, H. et al. 2023. Mixed-Type Tabular Data Synthesis with Score-based Diffusion in Latent Space.

[4] Solatorio, A. V. et al. 2023. REaLTabFormer: Generating Realistic Relational and Tabular Data using Transformers.

# Data Synthesis (How to evaluate?)



[1] DataCebo, Inc. Synthetic Data Metrics, 10 2023. Version 0.12.0.

[2] Xu, L. et al. 2019. Modeling tabular data using conditional gan.

[3] Zhao, Z. et al. 2021. Ctab-gan: Effective table data synthesizing.

# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection

## Generation



### **Purpose:**

Improve privacy while preserving the efficacy of LLMs



## Evaluation

### **Purpose:**

Evaluate the effectiveness and privacy risks of synthetic data in the age of LLMs.

# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection

## Generation



Propose a synthetic tabular data generation approach, which utilizes:

- I. KNN to maintain the efficacy
- II. fine-tuning LLMs for privacy preservation



## Evaluation

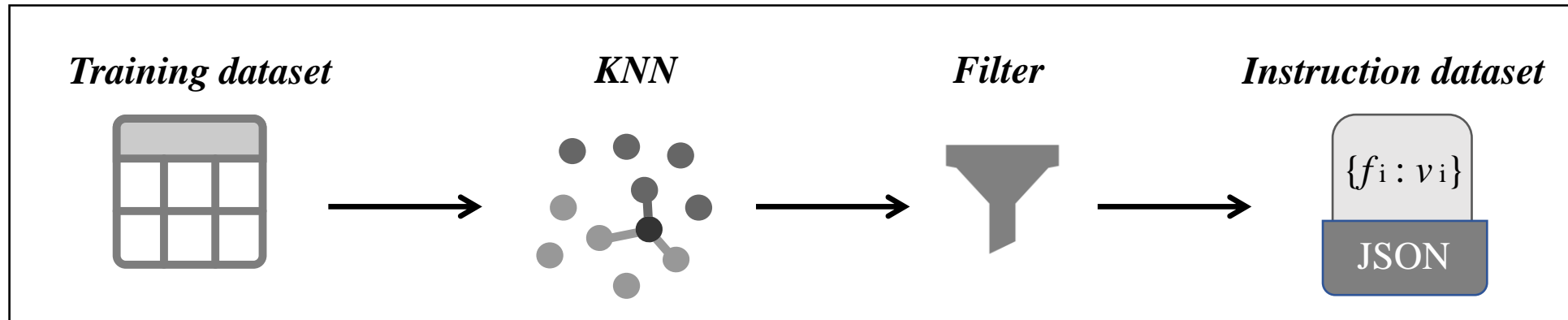
Introduce two new metrics to evaluate the **effectiveness** and **privacy** of synthetic data for LLM-based synthesis methods:

- I. LLM Effectiveness (LLE)
- II. Data Leakage Test (DLT)

# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection **Generation**

## ● I. Instruction dataset construction

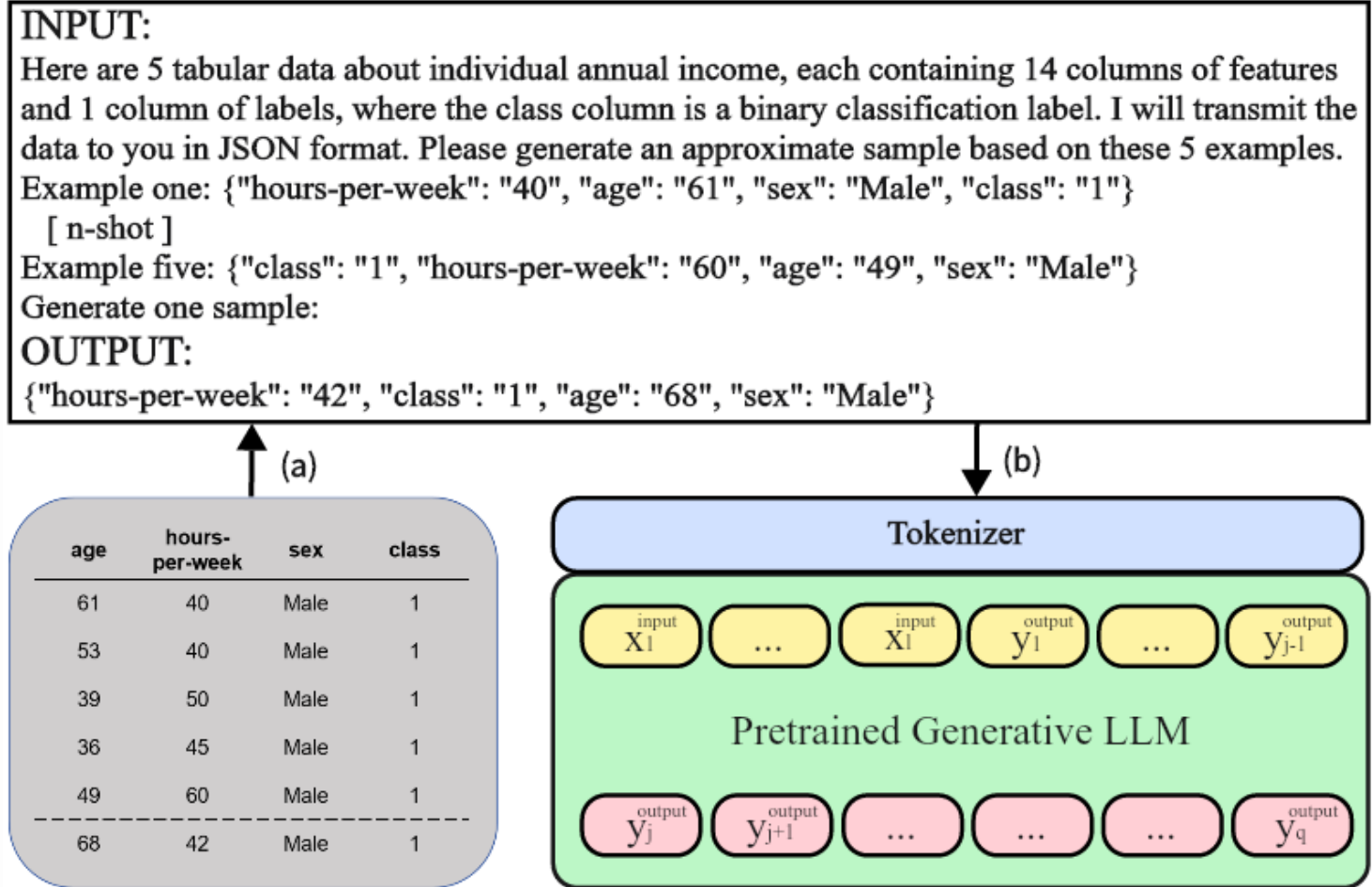
- Let the LLMs see the relationship between multiple similar rows and construct the structural tabular synthetic data format.
- To achieve this, we use the **KNN** algorithm to identify neighboring data for each instance.
- Every row of tabular data set gotten by KNN is converted into **JSON** dictionary format.



# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection

## Generation

- II. Instruction tuning based tabular data synthesizer
  - Fine-tune the LLM for the synthetic data generation task using the instruction dataset.
  - The **objective** of our fine-tuning strategy is to maximize the probability of generating the correct output sequence given the prompt describing the task and 5 input real data points.



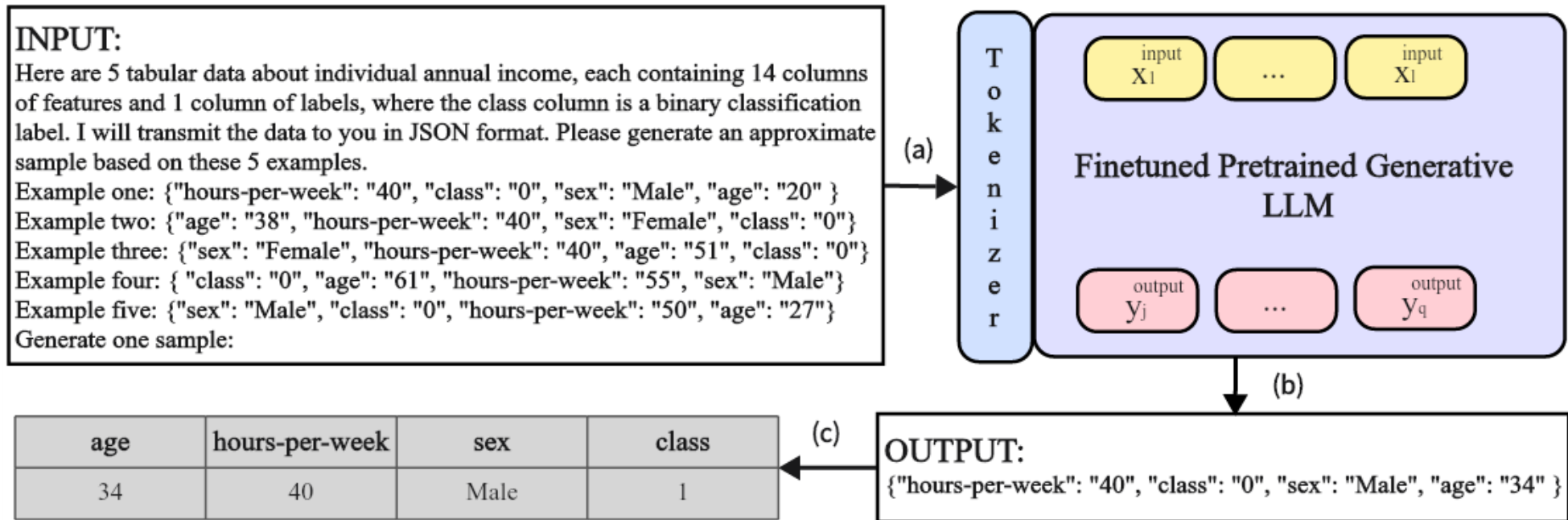


# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection

## Generation

### III. Sampling for synthetic data generation

- After the fine-tune, we denote the fine-tuned LLM as the generator for tabular data synthesis.
- The prompt dataset used for data generation is re-sampled randomly and is disjoint from the training set.
- Our method aims to teach the LLMs to extract patterns from the original data, unlike pretraining where the model primarily memorizes the data.



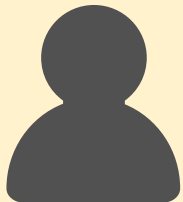
# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection

## Evaluation

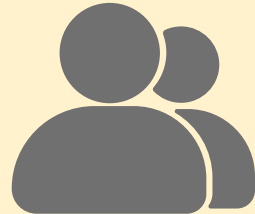
### ● Motivation

- With the development of LLMs, evaluating the effectiveness of synthetic data using weak classifiers(MLE) is losing its practical value and credibility.
- The commonly used data leakage metrics focus on measuring the "distance" between synthetic data and real data, without taking into account the extent to which the generator itself leaks data.

I trained a **SVM** classifier that works well on the diabetes dataset !



What? We no longer use SVM. We now use **CatBoost** or **LLMs**.



*original data*

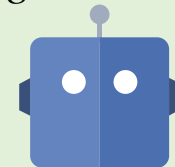


v.s.

*synthetic data*



*generator*



*original data*



# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection **Evaluation**

## LLE

- We propose using synthetic data to fine-tune a pretrained LLM and then evaluate the fine-tuned LLM on the real test set. We refer to this as LLM Effectiveness (LLE).
- LLMs trained on synthetic data should ideally achieve comparable or even surpass the performance of LLMs trained on real data, as measured by LLE on a real test set.

## DLT

- We introduce a new metric, Data Leakage Test (DLT), to quantify privacy by comparing the difference in perplexity (PPL) between a generator on real data and synthetic data as:  
$$DLT = PPL(D_{train}) - PPL(D_{syn})$$
- A higher DLT value indicates that the generator is more inclined to generate synthetic data rather than real data, thus reducing the likelihood of leaking real data.

# Experimental setup

- **Task & Data:** four classification datasets (German, Adult Income, Diabetes, Buddy)
- **Baselines:** Simple Augmentation (SMOTE), VAE-based (TVAE), GAN-based (CTABGAN), Diffusion-based (TabDDPM, TABSYN), LLM-based (REaLTabFormer, GReaT)
- **Evaluation Metric:** statistical characteristics (data\_mismatch, Wasserstein\_dist, CorrelationSimilarity), effectiveness (MLE, LLE), privacy (NewRowSynthesis, DCR, DLT)
- **Models:** we opt for LLaMA-2-7b-chat as the base model
- **Ablations:**
  - *w/o KNN:* Using randomly sampling instead of KNN in Instruction Dataset Construction
  - *w/o filter:* Don't use the filter operation in Instruction Dataset Construction
  - *w/o permutation:* Don't permute the features in Instruction Dataset Construction

[1] Hans Hofmann. Statlog (German Credit Data). UCI Machine Learning Repository, 1994.

[2] <https://github.com/vanderschaarlab/synthcity>

[3] Hugo Touvron, et al. Llama 2: Open foundation and fine-tuned chat models.

# Experiments – Main Results

Dataset	Metric	Original	HARMONIC	SMOTE	TVAE	CTAB	TabDDPM	TABSYN	RTF	GReaT
GM	DM	–	<b>0.00</b> $\pm$ 0.00	<u>0.14</u> $\pm$ 0.00	<u>0.14</u> $\pm$ 0.00	<u>0.14</u> $\pm$ 0.00	<u>0.14</u> $\pm$ 0.00	0.27 $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<u>0.14</u> $\pm$ 0.00
	WD	–	0.87 $\pm$ 0.07	0.85 $\pm$ 0.03	<u>0.70</u> $\pm$ 0.04	0.77 $\pm$ 0.02	0.73 $\pm$ 0.02	0.94 $\pm$ 0.09	<b>0.67</b> $\pm$ 0.01	0.93 $\pm$ 0.22
	CS	–	0.96	0.97	<b>0.99</b>	<u>0.98</u>	0.90	<u>0.98</u>	<u>0.98</u>	<u>0.98</u>
	MLE	0.50 $\pm$ 0.00	0.55 $\pm$ 0.03	<u>0.64</u> $\pm$ 0.02	0.61 $\pm$ 0.02	0.57 $\pm$ 0.02	<u>0.64</u> $\pm$ 0.01	0.63 $\pm$ 0.02	<b>0.65</b> $\pm$ 0.01	0.44 $\pm$ 0.03
	LLE	0.71 $\pm$ 0.00	0.64 $\pm$ 0.03	<u>0.67</u> $\pm$ 0.04	0.69 $\pm$ 0.03	0.71 $\pm$ 0.02	<u>0.67</u> $\pm$ 0.05	<b>0.72</b> $\pm$ 0.02	0.69 $\pm$ 0.03	0.55 $\pm$ 0.11
AD	DM	–	0.21 $\pm$ 0.15	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<u>0.06</u> $\pm$ 0.00
	WD	–	0.48 $\pm$ 0.15	0.49 $\pm$ 0.01	0.31 $\pm$ 0.04	0.07 $\pm$ 0.01	0.06 $\pm$ 0.01	0.07 $\pm$ 0.01	<b>0.03</b> $\pm$ 0.00	3.83 $\pm$ 0.09
	CS	–	0.90	<b>0.99</b>	<u>0.98</u>	0.97	<b>0.99</b>	0.97	<b>0.99</b>	0.94
	MLE	0.61 $\pm$ 0.00	0.67 $\pm$ 0.02	<u>0.75</u> $\pm$ 0.00	0.74 $\pm$ 0.00	0.73 $\pm$ 0.01	0.74 $\pm$ 0.00	0.73 $\pm$ 0.01	<b>0.76</b> $\pm$ 0.00	0.73 $\pm$ 0.01
	LLE	0.81 $\pm$ 0.00	0.80 $\pm$ 0.02	<u>0.84</u> $\pm$ 0.01	0.83 $\pm$ 0.01	0.83 $\pm$ 0.00	0.83 $\pm$ 0.00	0.81 $\pm$ 0.02	<b>0.85</b> $\pm$ 0.00	0.82 $\pm$ 0.02
DI	DM	–	<b>0.03</b> $\pm$ 0.05	<u>0.07</u> $\pm$ 0.05	0.10 $\pm$ 0.00	0.10 $\pm$ 0.00	<u>0.07</u> $\pm$ 0.05	0.10 $\pm$ 0.00	0.10 $\pm$ 0.00	<u>0.07</u> $\pm$ 0.05
	WD	–	0.14 $\pm$ 0.01	<b>0.07</b> $\pm$ 0.00	<b>0.07</b> $\pm$ 0.00	0.26 $\pm$ 0.01	<u>0.08</u> $\pm$ 0.00	0.09 $\pm$ 0.01	0.09 $\pm$ 0.02	0.13 $\pm$ 0.00
	CS	–	0.95	0.96	<u>0.98</u>	0.97	<b>0.99</b>	<u>0.98</u>	<b>0.99</b>	0.88
	MLE	0.56 $\pm$ 0.00	0.46 $\pm$ 0.02	<b>0.72</b> $\pm$ 0.03	<u>0.71</u> $\pm$ 0.02	0.67 $\pm$ 0.02	<u>0.71</u> $\pm$ 0.02	0.68 $\pm$ 0.03	0.66 $\pm$ 0.03	0.45 $\pm$ 0.03
	LLE	0.70 $\pm$ 0.00	<u>0.75</u> $\pm$ 0.00	0.69 $\pm$ 0.04	<u>0.72</u> $\pm$ 0.04	0.62 $\pm$ 0.09	<u>0.72</u> $\pm$ 0.03	<b>0.77</b> $\pm$ 0.01	0.70 $\pm$ 0.04	0.71 $\pm$ 0.03
BU	DM	–	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<b>0.00</b> $\pm$ 0.00	<u>0.03</u> $\pm$ 0.04
	WD	–	0.48 $\pm$ 0.16	0.23 $\pm$ 0.02	0.10 $\pm$ 0.01	<u>0.05</u> $\pm$ 0.00	0.06 $\pm$ 0.02	<b>0.04</b> $\pm$ 0.00	<b>0.04</b> $\pm$ 0.00	2292.47 $\pm$ 1014.22
	CS	–	0.93	0.98	0.97	<u>0.99</u>	<b>1.00</b>	<b>1.00</b>	0.98	<b>1.00</b>
	MLE	0.38 $\pm$ 0.00	<b>0.27</b> $\pm$ 0.03	0.25 $\pm$ 0.02	<b>0.27</b> $\pm$ 0.03	0.26 $\pm$ 0.01	<b>0.27</b> $\pm$ 0.01	0.26 $\pm$ 0.01	<u>0.26</u> $\pm$ 0.00	0.24 $\pm$ 0.03
	LLE	0.88 $\pm$ 0.00	0.82 $\pm$ 0.03	0.85 $\pm$ 0.04	<b>0.86</b> $\pm$ 0.01	0.82 $\pm$ 0.02	<u>0.85</u> $\pm$ 0.01	<b>0.86</b> $\pm$ 0.01	<u>0.70</u> $\pm$ 0.14	0.81 $\pm$ 0.03

- Our method exhibits performance comparable to existing SOTA approaches, especially from a statistical perspective, our method approaches optimality.
- In the era of LLMs, detecting synthetic data on LLMs is essential but it doesn't necessarily have to align with the performance of machine learning models.

# Experiments – Main Results

Dataset	Metric	HARMONIC	SMOTE	TVAE	CTAB	TabDDPM	TABSYN	RTF	GReaT
GM	NRS	<b>1.00</b>	1.00	1.00	1.00	1.00	1.00	1.00	1.00
	DCR	<b>8.08</b>	2.77	4.09	5.36	2.21	3.98	4.60	<u>5.84</u>
	DLT	<b>-0.16</b>	—	—	—	—	—	-22.04	<u>-2.14</u>
AD	NRS	<b>1.00</b>	0.95	1.00	1.00	1.00	1.00	1.00	1.00
	DCR	<b>2.47</b>	0.16	0.49	0.82	0.50	0.86	0.57	<u>1.51</u>
	DLT	<u>-0.98</u>	—	—	—	—	—	-163.71	<b>-0.67</b>
DI	NRS	<b>1.00</b>	1.00	1.00	1.00	1.00	1.00	1.00	1.00
	DCR	0.44	0.28	0.33	0.72	0.21	<b>1.37</b>	0.36	<u>1.36</u>
	DLT	<b>-0.37</b>	—	—	—	—	—	-42.46	<u>-0.44</u>
BU	NRS	<b>1.00</b>	0.93	1.00	1.00	0.99	1.00	1.00	1.00
	DCR	<u>2.52</u>	0.15	0.66	0.70	0.18	1.38	0.38	<b>8.30</b>
	DLT	<b>-0.34</b>	—	—	—	—	—	-41.13	<u>-2.22</u>

- Our method prioritizes privacy in the synthetic data generation, HARMONIC surpasses or comes in a close second for all datasets across all metrics.
- DLT indicator suggests that LLMs pose a significant risk of data leakage.



# Conclusion

## Contributions

---

- We recognize that it is crucial to not only focus on the strong data generation ability of LLM in this era, but also pay attention to the potential privacy risks it may bring.
- We develop a framework, HARMONIC, for tabular data synthesis based on LLM. Our method ensures privacy preservation while maintaining the effectiveness of synthetic data.
- Under the HARMONIC framework, a set of metrics is proposed for the effectiveness in downstream LLMs tasks and privacy risk evaluation of synthetic tabular data.

## Future Work

---

- We will explore tasks beyond categorical datasets.
- We will also investigate data formats beyond tabular data.
- We will attempt to implement our methods in real-world applications.

*Refer to our paper for more experiments and discussions*

# HARMONIC: Harnessing LLMs for Tabular Data Synthesis and Privacy Protection

**Yuxin Wang, Duanyu Feng, Yongfu Dai, Zhengyu Chen, Jimin Huang,  
Sophia Ananiadou, Qianqian Xie, Hao Wang**



四川大學  
SICHUAN UNIVERSITY

