# Poster: NETDPSYN: Synthesizing Network Traces under Differential Privacy

Danyu Sun[†], Joann Qiongna Chen[§], Chen Gong[‡], Tianhao Wang[‡], and Zhou Li[†✉]

[†]University of California, Irvine, [§]San Diego State University, [‡]The University of Virginia

{danyus2, zhou.li}@uci.edu, {jchen27}@sdsu.edu, {fzv6en, tianhao}@virginia.edu,

## Abstract

As the utilization of network traces for the network measurement research becomes increasingly prevalent, concerns regarding privacy leakage from network traces have garnered the public's attention. To safeguard network traces, researchers have proposed the *trace synthesis* that retains the essential properties of the raw data. However, previous works also show that synthesis traces with generative models are vulnerable under linkage attacks.

This paper introduces NETDPSYN, the first system to synthesize high-fidelity network traces under privacy guarantees. NETDPSYN is built with the Differential Privacy (DP) framework as its core, which is significantly different from prior works that apply DP when training the generative model. The experiments conducted on three flow and two packet datasets indicate that NETDPSYN achieves much better data utility in downstream tasks like anomaly detection. NETDPSYN is also 2.5 times faster than the other methods on average in data synthesis.

## ACKNOWLEDGMENT

## REFERENCES

[1] D. Sun, J. Q. Chen, C. Gong, T. Wang, and Z. Li. Netdpsyn: Synthesizing network traces under differential privacy. In *Proceedings of the 2024 ACM on Internet Measurement Conference*, IMC '24, page 545–554, New York, NY, USA, 2024. Association for Computing Machinery.

---

[✉] Corresponding authors.
[1]https://dl.acm.org/doi/10.1145/3646547.3689011.
[2]https://arxiv.org/pdf/2409.05249.

Danyu Sun[1]   Joann Qiongna Chen[2]   Chen Gong[3]   Tianhao Wang[3]   Zhou Li[1]

[1]University of California, Irvine    [2]San Diego State University    [3]University of Virginia
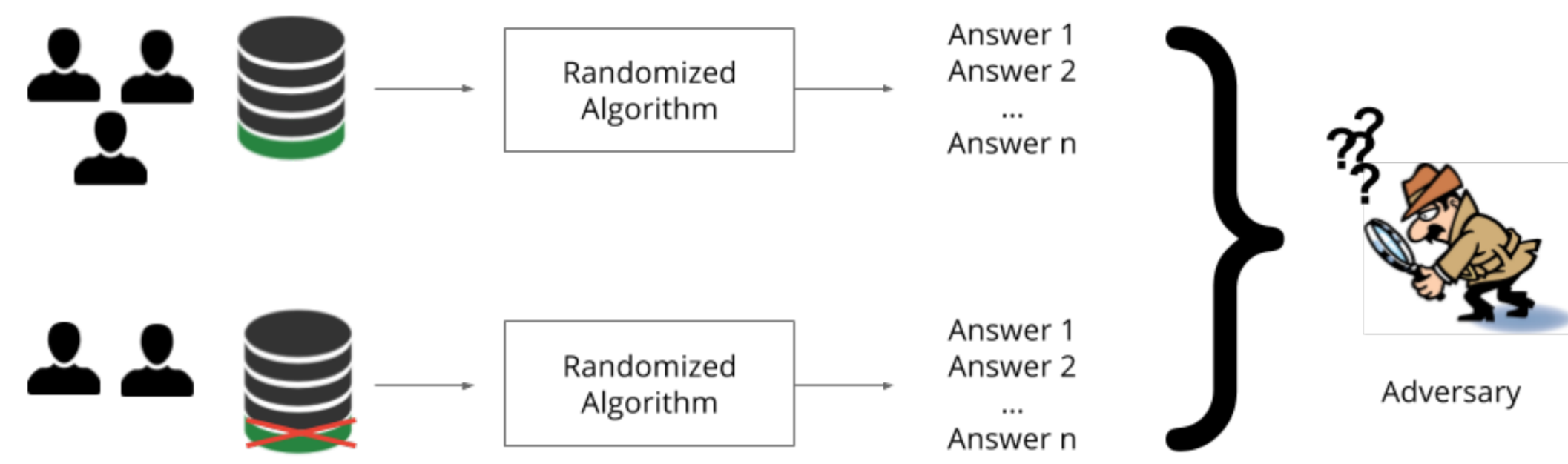
## Background

- **Network Data and Privacy:** We consider header fields of network *packet* or *flows* as the target for data synthesis. Releasing the header without payload still raises privacy concerns, and the main solutions include data anonymization and synthesis.
- **Differential Privacy:** A randomized mechanism $\mathcal{A}$ satisfies $(\varepsilon, \delta)$-differential privacy ($\varepsilon > 0$ and $\delta > 0$), if and only if, for any two neighboring datasets $D$ and $D'$, it holds that,
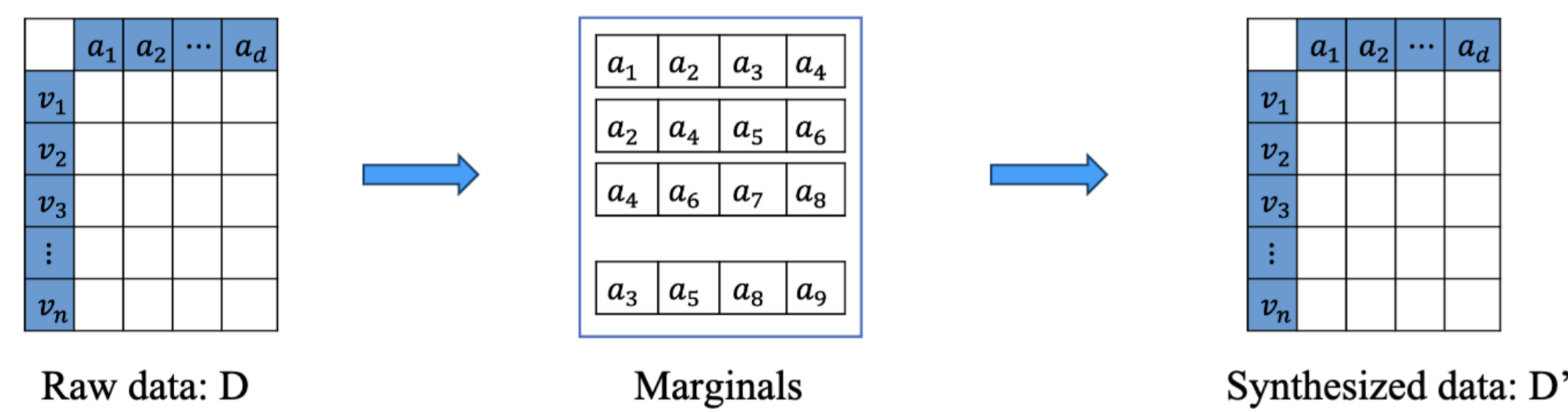
$$\Pr[\mathcal{A}(D) \in \mathcal{O}] \leq e^{\varepsilon} \Pr[\mathcal{A}(D') \in \mathcal{O}] + \delta \qquad (1)$$

DP ensures the privacy of an individual's data within a dataset is preserved data processing. It guarantees that the result of any computation, such as a database query, remains essentially unchanged whether or not any single individual's data is included or excluded.
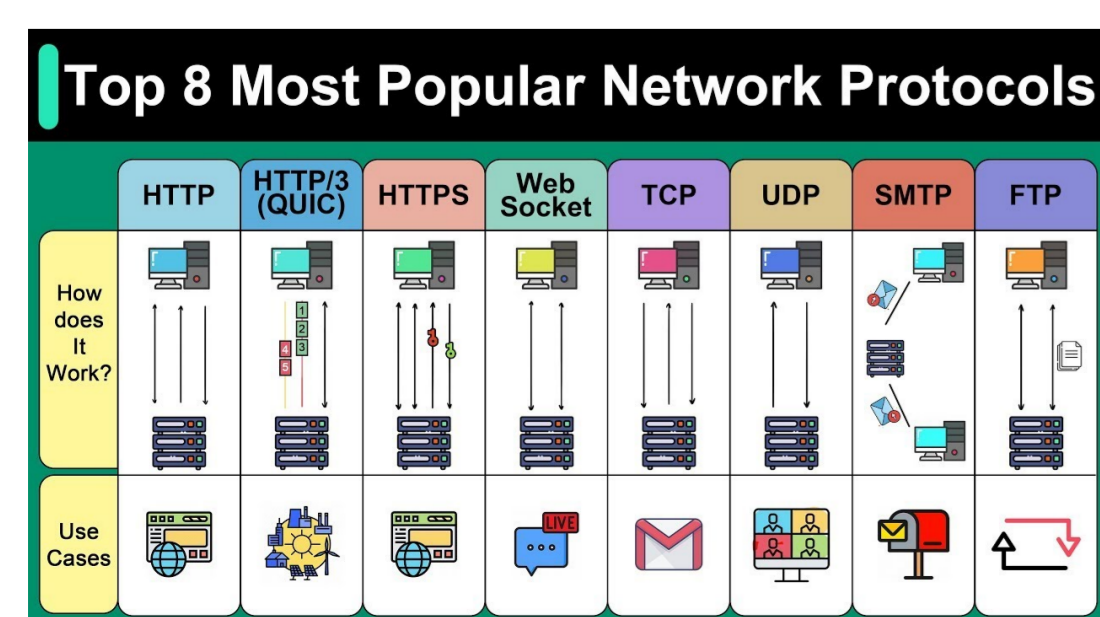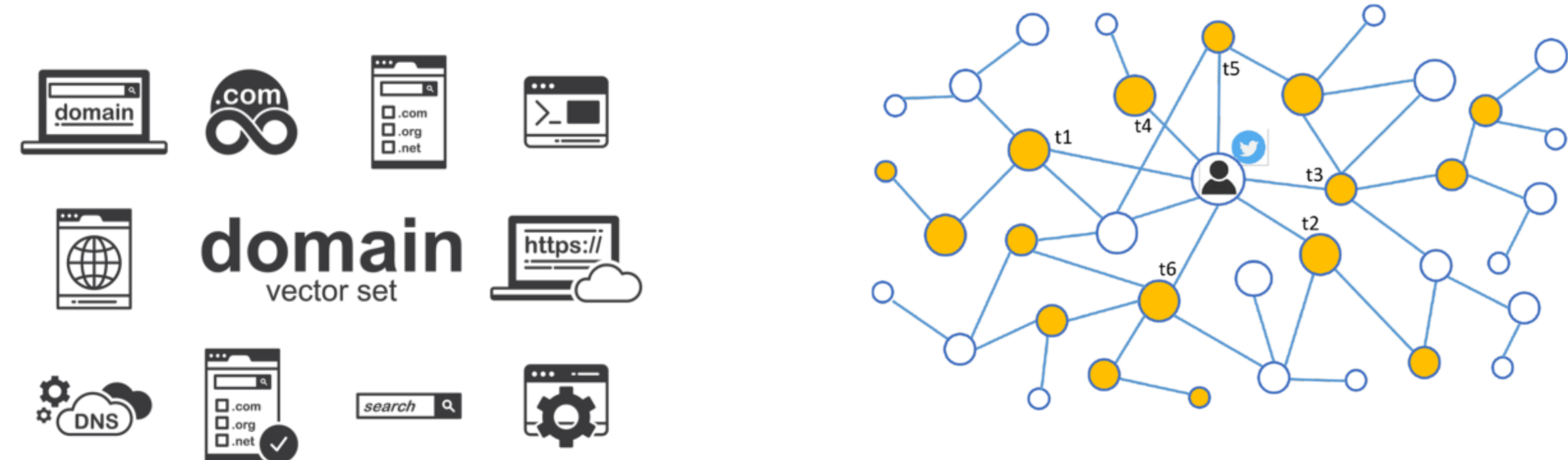


## Motivation

- To achieve formal privacy guarantees, the training of the generative model can be hardened with DP. We found **Only two** works that synthesize network traces applied DP [1, 3]. However, the **data utility** will be significantly worse.
- We pursue a different direction to capture the **underlying distribution** of the original data and then synthesizing network records **after** they are protected by DP.
- PrivSyn [4] handles high-dimensional datasets by automatically selecting and constructing noisy marginal tables that captures the data distribution.



Raw data: D          Marginals          Synthesized data: D'
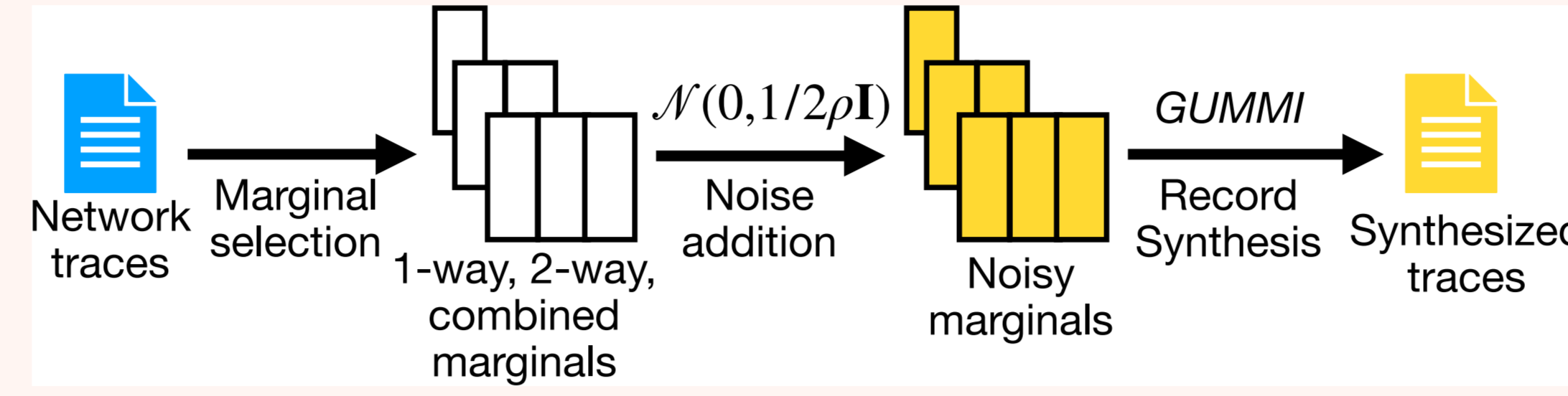
## Challenges

- **Large Domain Size:** E.g., the number of bytes per flow and long-tail distribution.
- **Temporal Patterns:** It is essential to downstream applications. However, such patterns cannot be synthesized using marginals tables.
- **Inherent Constraints of Network Protocol:** E.g., FTP communications use TCP and the destination ports are 20 and 21.
- **Time Consuming:** PrivSyn [4] uses the Gradually Update Method (GUM) to synthesize records from noisy marginal tables, but this method can take a long time to converge for large network datasets.



**Top 8 Most Popular Network Protocols**

TIME-CONSUMING

## NetDPSyn [2]

- **Input**: Private Network Traces.
- **Output**: Synthesized Data.



## Pre-Process

- **IP:** Bin the low-count IP address by the /30 prefix.
- **Port:** Keep common port under 1024, bin higher by 10.
- **Integer and float-point:** Log transformation.
- **Timestamp:** Group-wise differences between timestamp.
- **Category:** No need to be binned.

## Marginal Selection

- **Marginal Tables:** Marginal tables example for `dstport` and `type` computed on TON dataset.

| $v$ | $M_d(v)$ |
|---|---|
| $\langle 53, * \rangle$ | 82828 |
| $\langle 80, * \rangle$ | 68748 |
| $\langle 15600, * \rangle$ | 27255 |

**(a) 1-way marginal for `dstport`.**

| $v$ | $M_t(v)$ |
|---|---|
| $\langle *, \text{normal} \rangle$ | 166494 |
| $\langle *, \text{injection} \rangle$ | 15951 |

**(b) 1-way marginal for `type`.**

| $v$ | $M_{dt}(v)$ |
|---|---|
| $\langle 53, \text{normal} \rangle$ | 74547.08 |
| $\langle 53, \text{injection} \rangle$ | 554.71 |
| $\langle 80, \text{normal} \rangle$ | 12297.88 |
| $\langle 80, \text{injection} \rangle$ | 15396.66 |
| $\langle 15600, \text{normal} \rangle$ | 27247.02 |
| $\langle 15600, \text{injection} \rangle$ | 20.09 |

**(c) 2-way noisy marginal before marginal post-processing.**

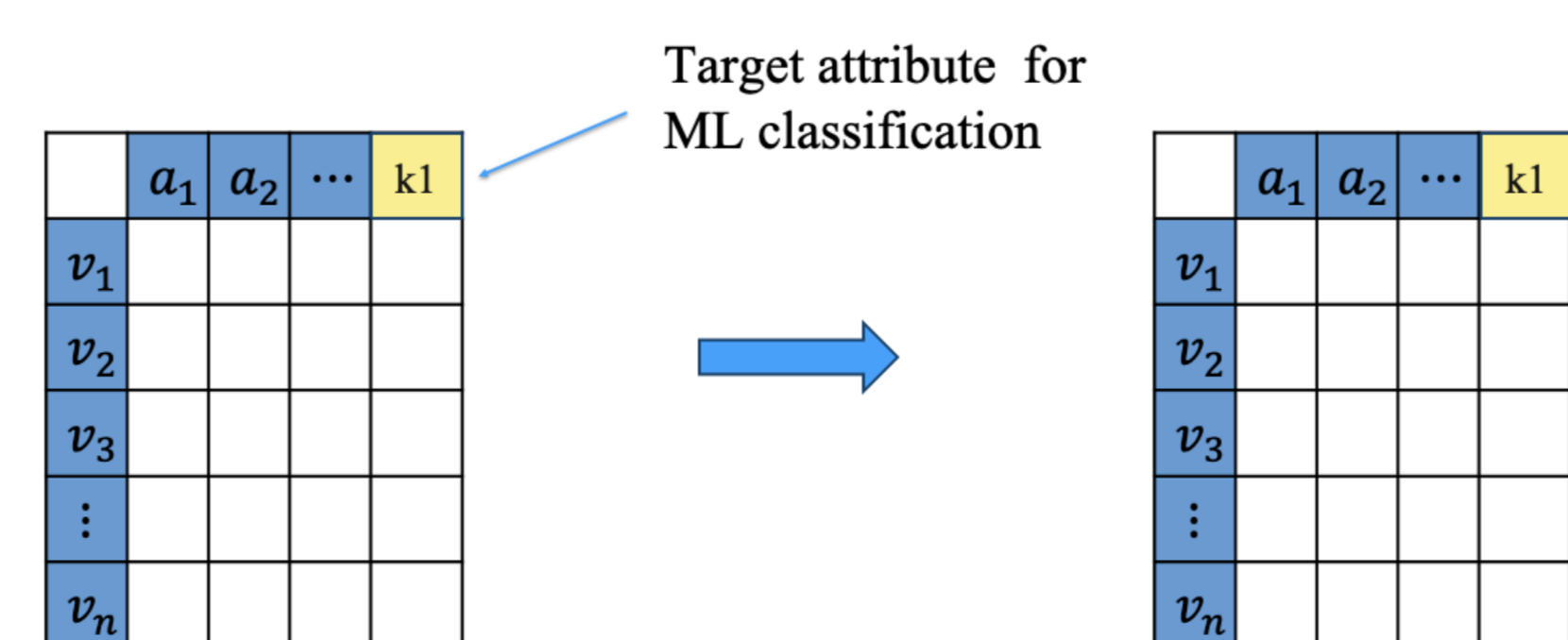| $v$ | $\tilde{M}_{dt}(v)$ |
|---|---|
| $\langle 53, \text{normal} \rangle$ | 74566 |
| $\langle 53, \text{injection} \rangle$ | 558 |
| $\langle 80, \text{normal} \rangle$ | 12308 |
| $\langle 80, \text{injection} \rangle$ | 15364 |
| $\langle 15600, \text{normal} \rangle$ | 27255 |
| $\langle 15600, \text{injection} \rangle$ | 0 |

**(d) Actual 2-way marginal.**

- **Marginal Selection:** DenseMarg [4] formalizes the marginal selection problem as an optimization problem that balances dependency error (error caused by missing a marginal) and noise error (error caused by adding noises to a selected marginal).

*noise error*     *dependency error*

$$\text{minimize} \sum_{i=1}^{m} [\psi_i x_i + \phi_i(1 - x_i)] \text{ s.t. } x_i \in \{0, 1\}$$

*selected two-way marginal*     *non-selected two-way marginal*

## GUMMI



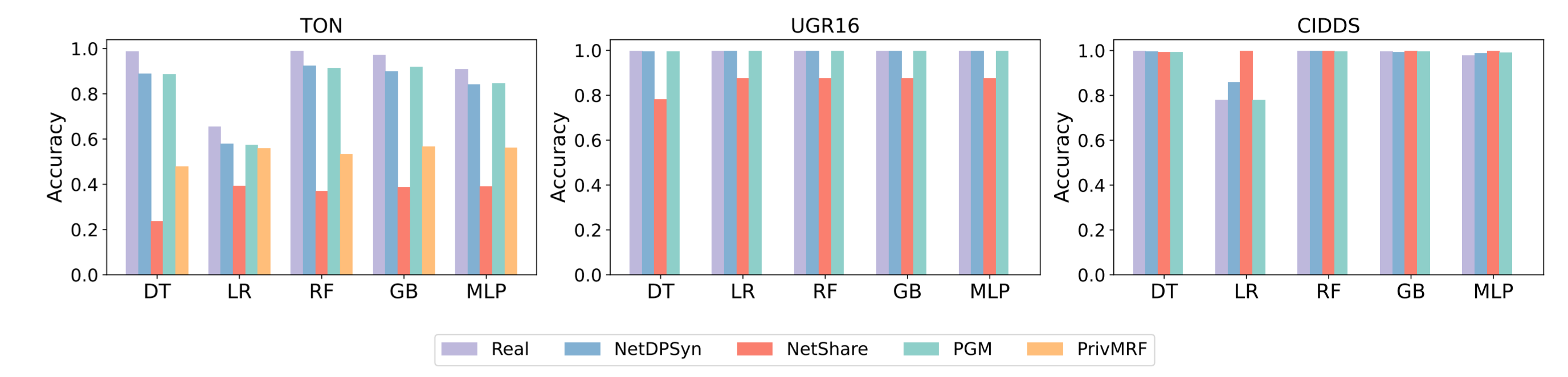Target attribute for ML classification

Step 1: Initialize dataset that **contains key** to the downstream tasks, e.g., label

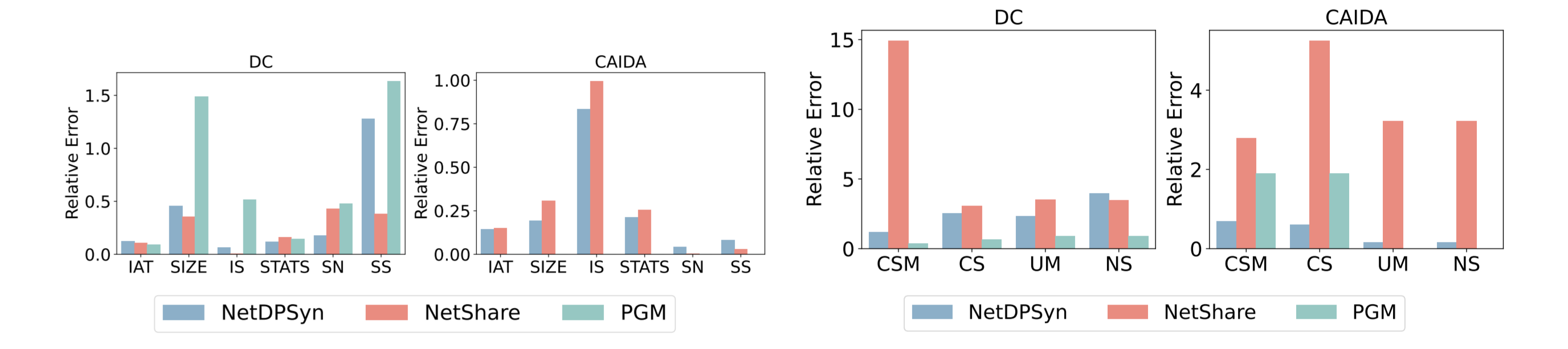Step 2: Using all noisy marginal to update the initialized dataset.

## Evaluation Results

- **Machine Learning Classification Tasks on 3 Flow Data:**



Accuracy of 5 ML Tasks.

- **Anomaly Detection and Sketch Algorithms on 2 Packet Data:**



Relative Error of NetML.          Relative Error of Various Sketch Algorithms.

- **Running Time of Each Method in Minutes:**

|  | NetDPSyn | NetShare | PGM | PrivMRF |
|---|---|---|---|---|
| TON | 10 | 27 | 70 | 240 |
| CIDDS | 20 | 100 | 55 | N/A |
| UGR16 | 40 | 94 | 55 | N/A |
| CAIDA | 35 | 30 | 54 | N/A |
| DC | 20 | 100 | 24 | N/A |

## Conclusions and Future Work

Conclusions:

- Synthesize **high-fidelity** network traces under privacy guarantee.
- Achieve **better data utility** in downstream task.
- **2.5 times** faster than other methods.

Future Work:

- Model the temporal pattern in complex representation.
- Cover all types of network environment and data type.
- Evaluate advanced downstream tasks like graph-based anomaly detection.

## Project Resources



Paper Link          GitHub Repository          Danyu Sun's Homepage

## References

[1] Liyue Fan and Akarsh Pokkunuru.
Dpnet: Differentially private network traffic synthesis with generative adversarial networks.
In *IFIP Annual Conference on Data and Applications Security and Privacy*, pages 3–21. Springer, 2021.

[2] Danyu Sun, Joann Qiongna Chen, Chen Gong, Tianhao Wang, and Zhou Li.
Netdpsyn: Synthesizing network traces under differential privacy.
In *Proceedings of the 2024 ACM on Internet Measurement Conference*, IMC '24, page 545–554, New York, NY, USA, 2024. Association for Computing Machinery.

[3] Yucheng Yin, Zinan Lin, Minhao Jin, Giulia Fanti, and Vyas Sekar.
Practical gan-based synthetic ip header trace generation using netshare.
In *Proceedings of the ACM SIGCOMM 2022 Conference*, pages 458–472, 2022.

[4] Zhikun Zhang, Tianhao Wang, Ninghui Li, Jean Honorio, Michael Backes, Shibo He, Jiming Chen, and Yang Zhang.
{PrivSyn}: Differentially private data synthesis.
In *30th USENIX Security Symposium (USENIX Security 21)*, pages 929–946, 2021.