



university of
groningen

faculty of science
and engineering

bernoulli institute



INFORMATION
SYSTEMS GROUP
ENABLING FLEXIBILITY

Federated Synthetic Data Generation with Stronger Security Guarantees

Ali Reza Ghavamipour

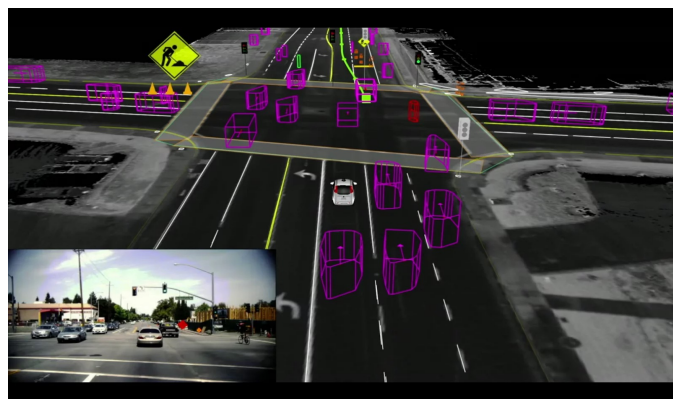
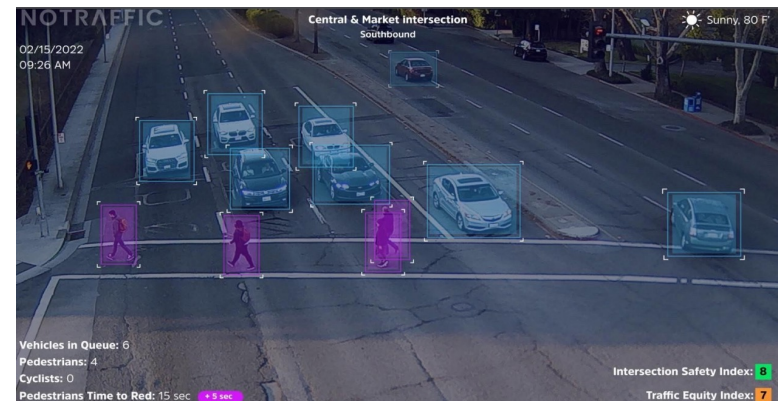
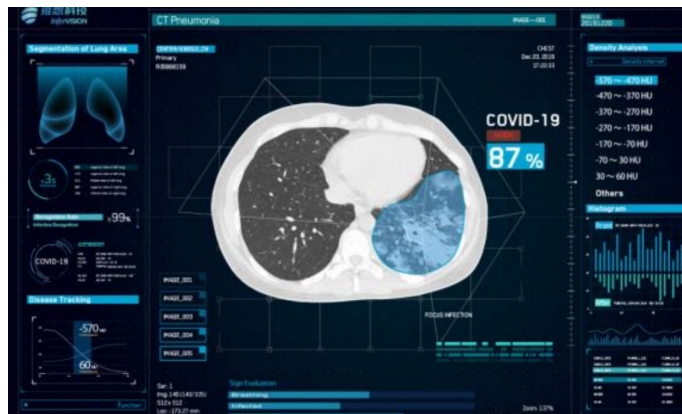
Rui Wang

Fatih Turkmen

Kaitai Liang



Machine Learning





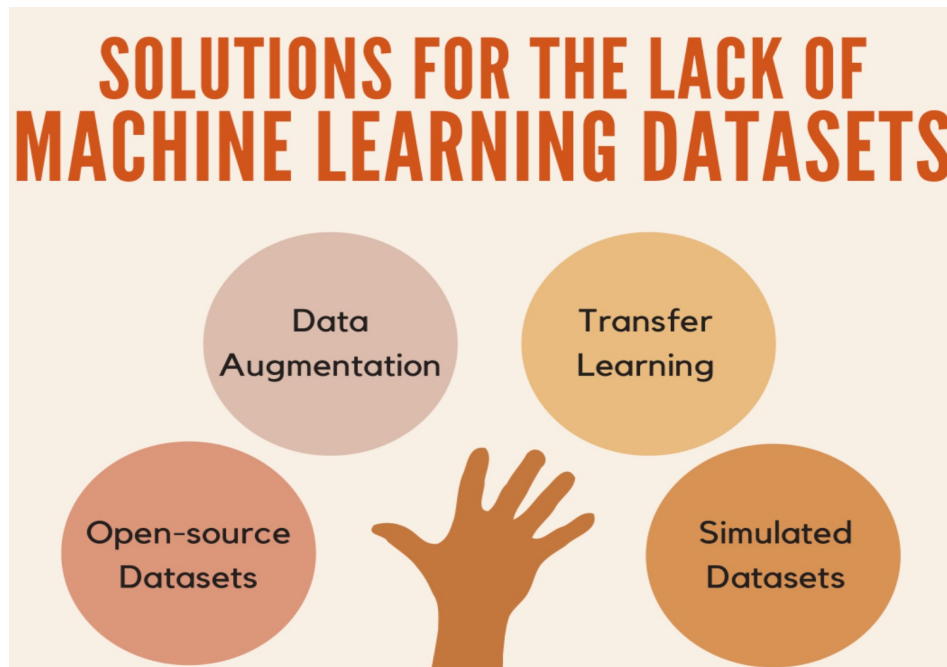
Data ...

- › More data = More accurate model

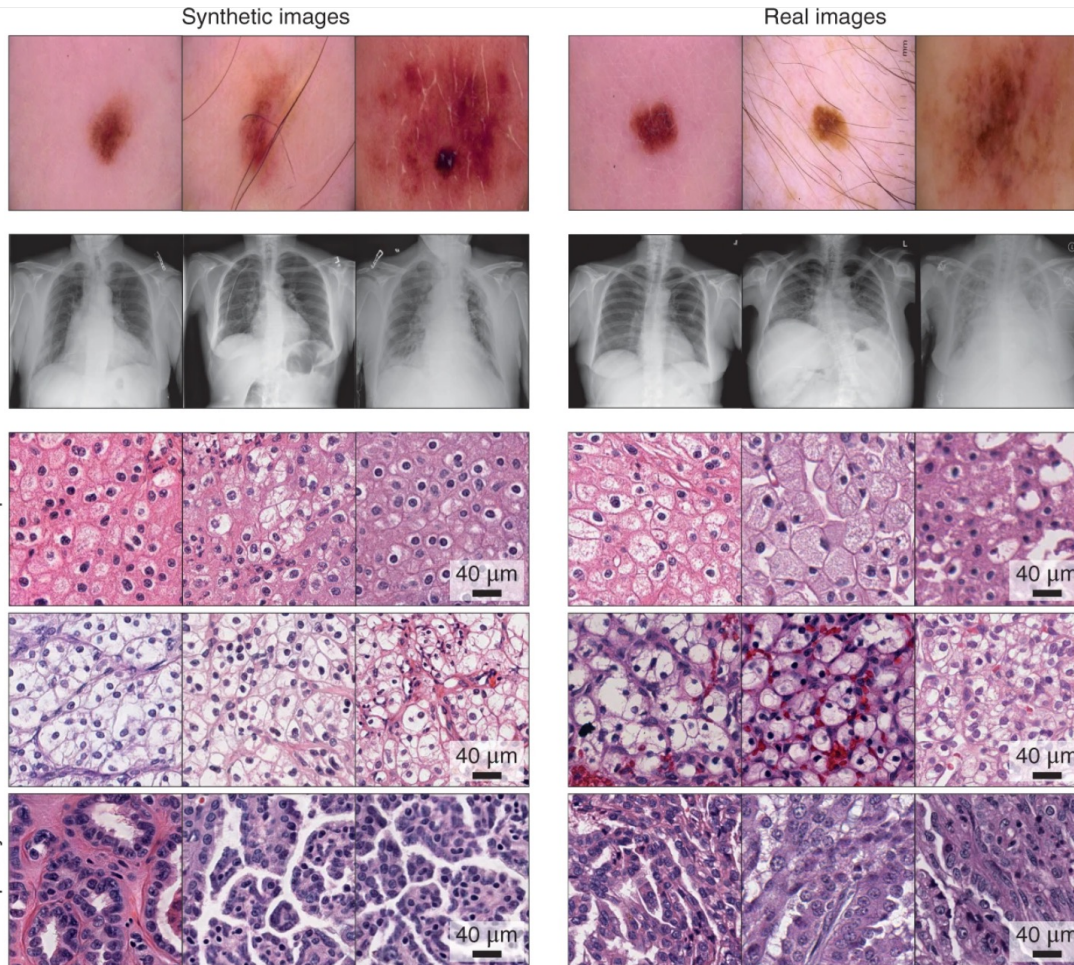
- › Restrictions on collecting real data:
 1. Regulations
 2. Privacy concerns
 3. Competitive advantage
 4. Liability



Possible solutions?



+ Synthetic data!



<https://www.nature.com/articles/s41551-021-00751-8>



Generative Adversarial Networks (GAN)

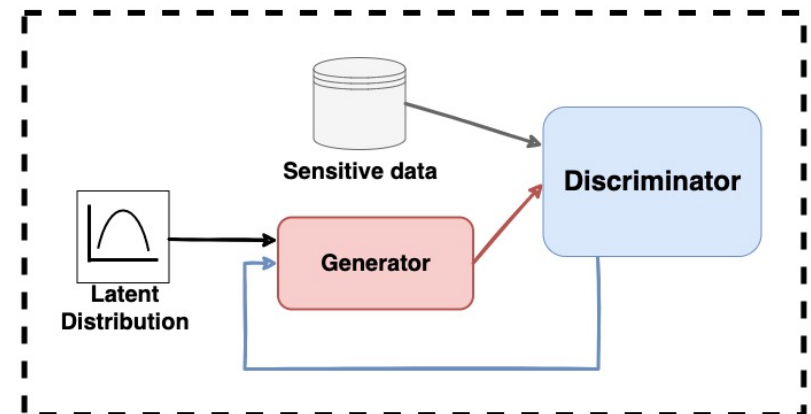
A class of machine learning frameworks
 by Ian Goodfellow in 2014.

Two deep neural networks, generator (G)
 and discriminator (D).

Zero-sum game between discriminator
 and generator.

Generator learns to produce new data with
 the same statistics as the training set.

Computationally expensive.





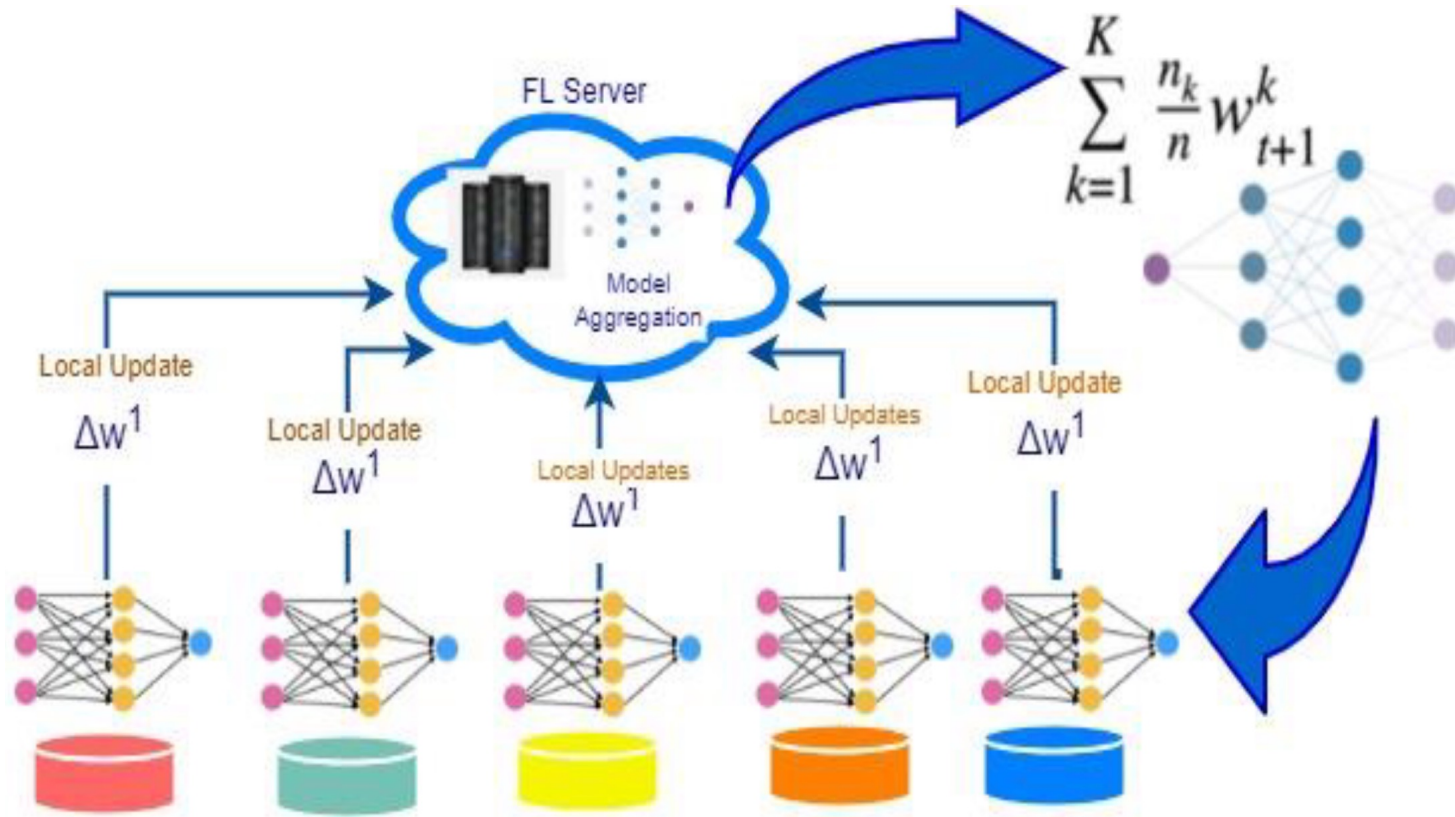
Do we have enough data?

- › But in the real world data is often distributed across multiple data holders.

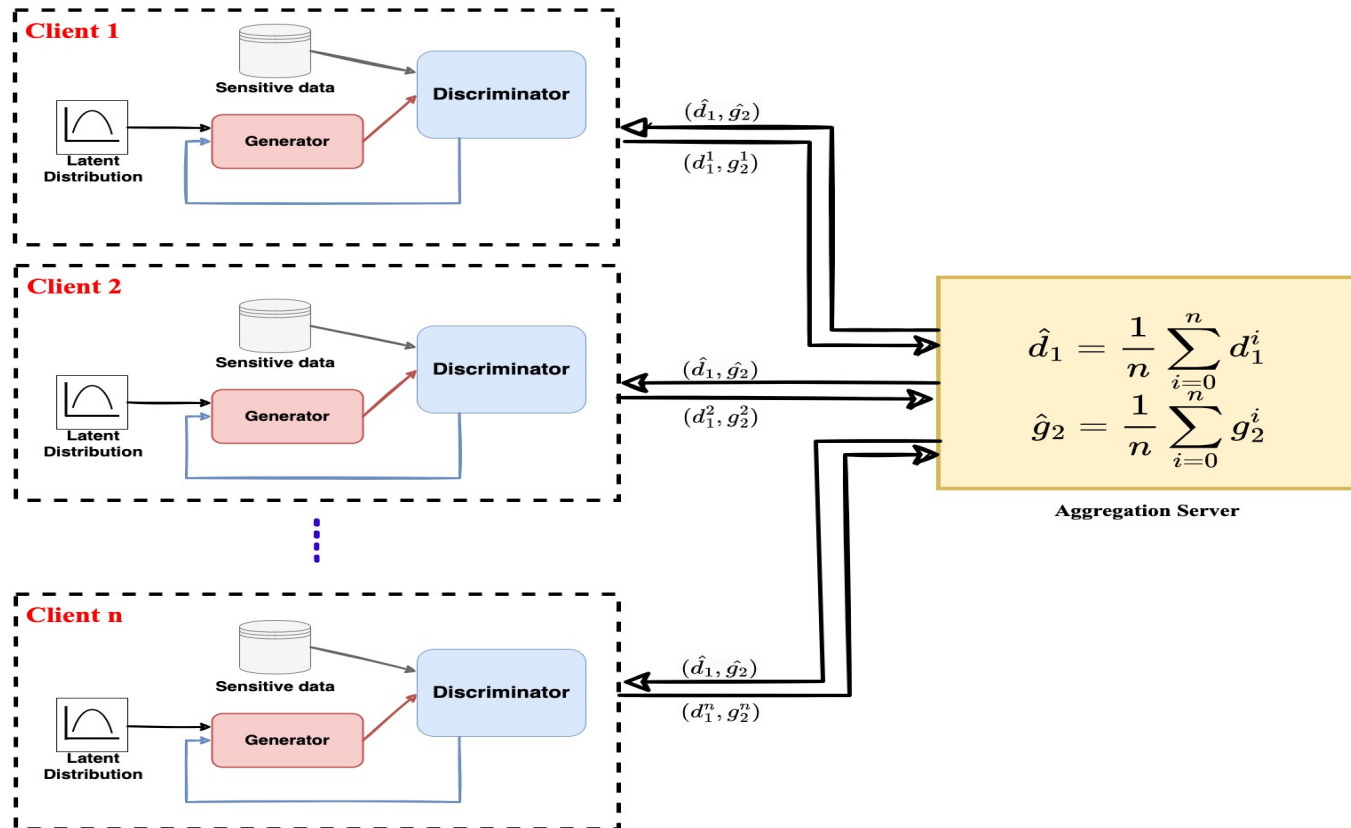
Privacy?



Federated Learning

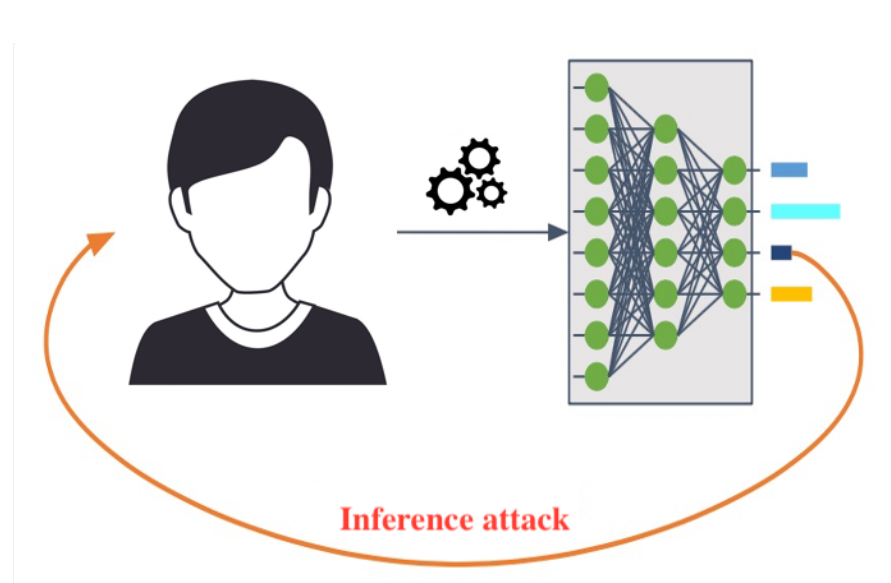


Federated GAN



Privacy Attacks Against ML and FL

- › Reconstruction Attacks
- › Membership Inference Attacks
- › Property Inference Attacks





Threat Model

1. Aggregation Server
2. Client (data owners)

Honest-but-curious:

- › Follow the protocol and not deviate from the protocol steps
- › May extract information about private data by analyzing the output or intermediate results



Possible Solutions

- › Encryption-based technique:
 - Homomorphic Encryption
 - Multi-Party Computation

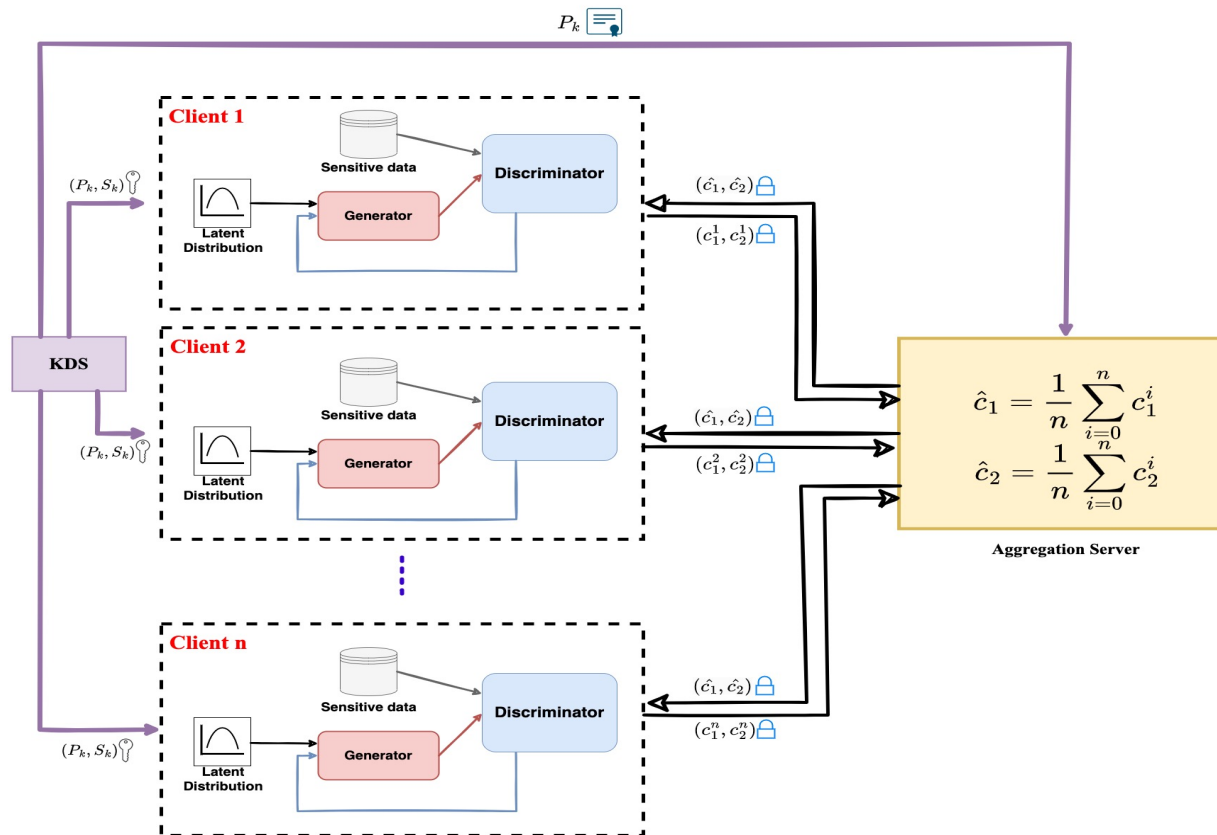
- › Perturbation-based techniques:
 - Differential Privacy



Homomorphic Encryption (HE)

- › Allows users to perform computations over encrypted data without decrypting it.
- › The result of computation remains in encrypted form.
- › We use CKKS scheme in this work.
- › CKKS uses polynomial approximation
- › Ability to perform computation on large-scale continuous data, such as matrices.
- › CKKS is Leveled HE scheme that supports a fixed number of "multiplication depth" levels.

Federated GAN + HE





All the issues solved?

- › Can homomorphic encryption also protect the sensitive data against semi-honest **client**?
- › No!
- › Solution?
- › **Differential Privacy.**



Differential Privacy (DP)

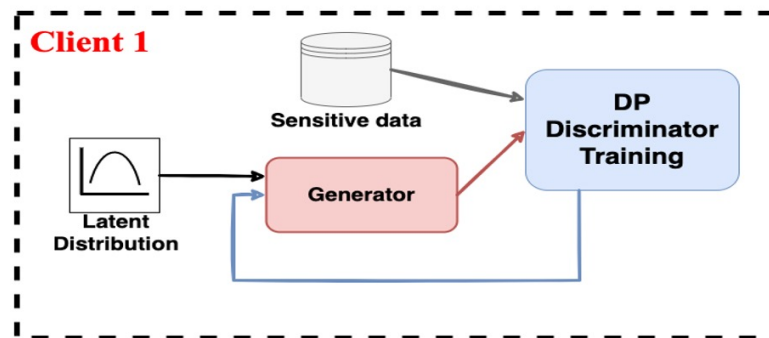
- › A mathematical framework for protecting the privacy by adding random noise
- › Enables data analysis while protecting the privacy of sensitive information

A process A is ϵ -differentially private if for all databases D_1 and D_2 which differ in only one individual:

$$\mathbb{P}[A(D_1) = O] \leq e^\epsilon \cdot \mathbb{P}[A(D_2) = O]$$

Differential Privacy (DP)

- › Adding Gaussian noise **during** the model training
- › Only to discriminator



$$\mathbf{g}^{(t)} := \nabla_{\theta} \mathcal{L}(\theta, \mathbf{w})$$

(Compute the per-sample gradients)

$$\hat{\mathbf{g}}^{(t)} := \mathcal{A}_{\sigma, C}(\mathbf{g}^{(t)}) = \text{clip}(\mathbf{g}^{(t)}, C) + \mathcal{N}(0, \sigma^2 C^2 I)$$

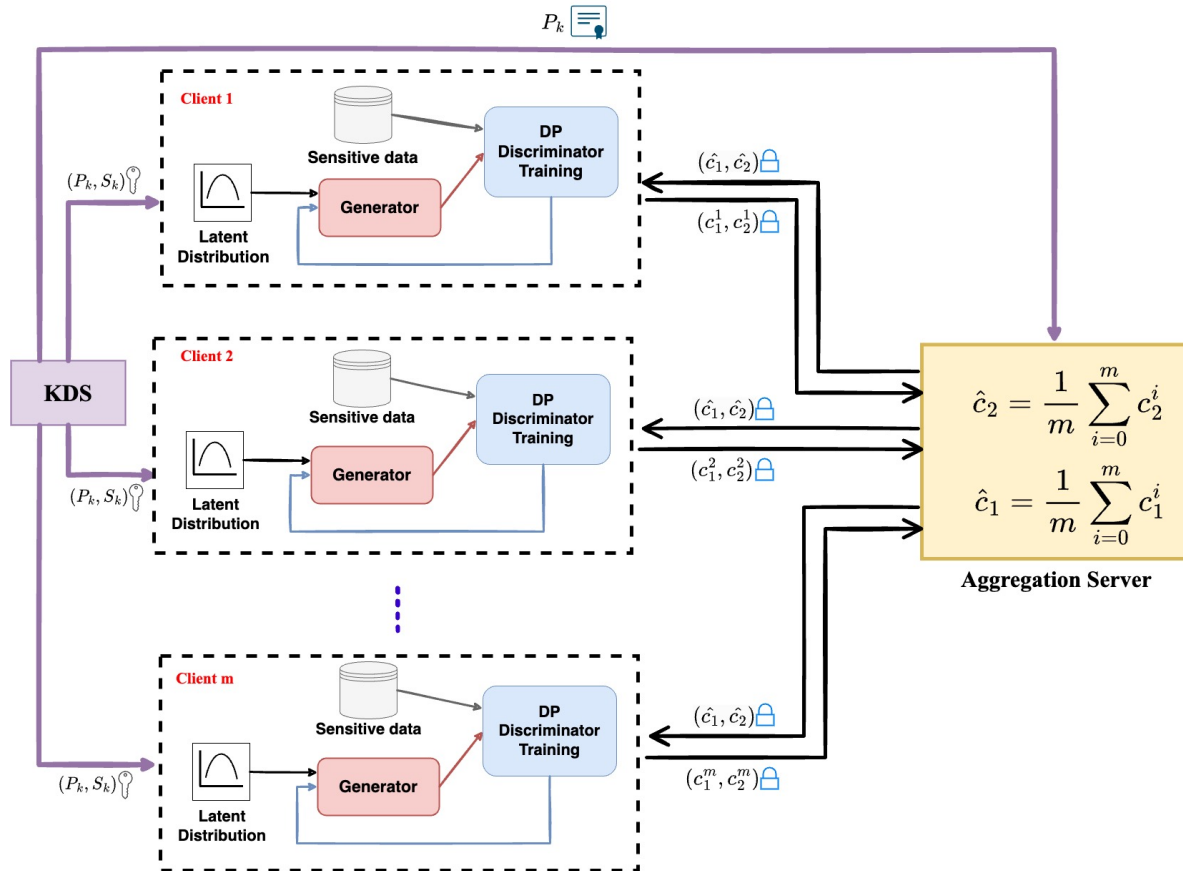
(Clipping and noise addition)

$$\theta^{(t+1)} := \theta^{(t)} - \eta_D \cdot \hat{\mathbf{g}}^{(t)}$$

(Gradient descent step)



Federated GAN + HE + DP = (PP-FedGAN)





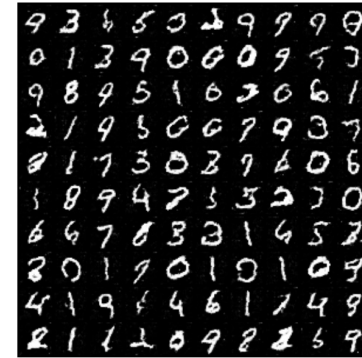
Experiments



(a) $\epsilon = 8.49$



(b) $\epsilon = 2.47$



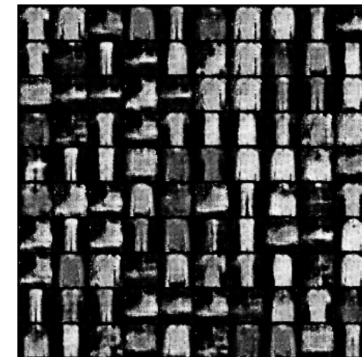
(c) $\epsilon = 1.39$



(d) $\epsilon = 6.2$



(e) $\epsilon = 2.1$



(f) $\epsilon = 1.3$



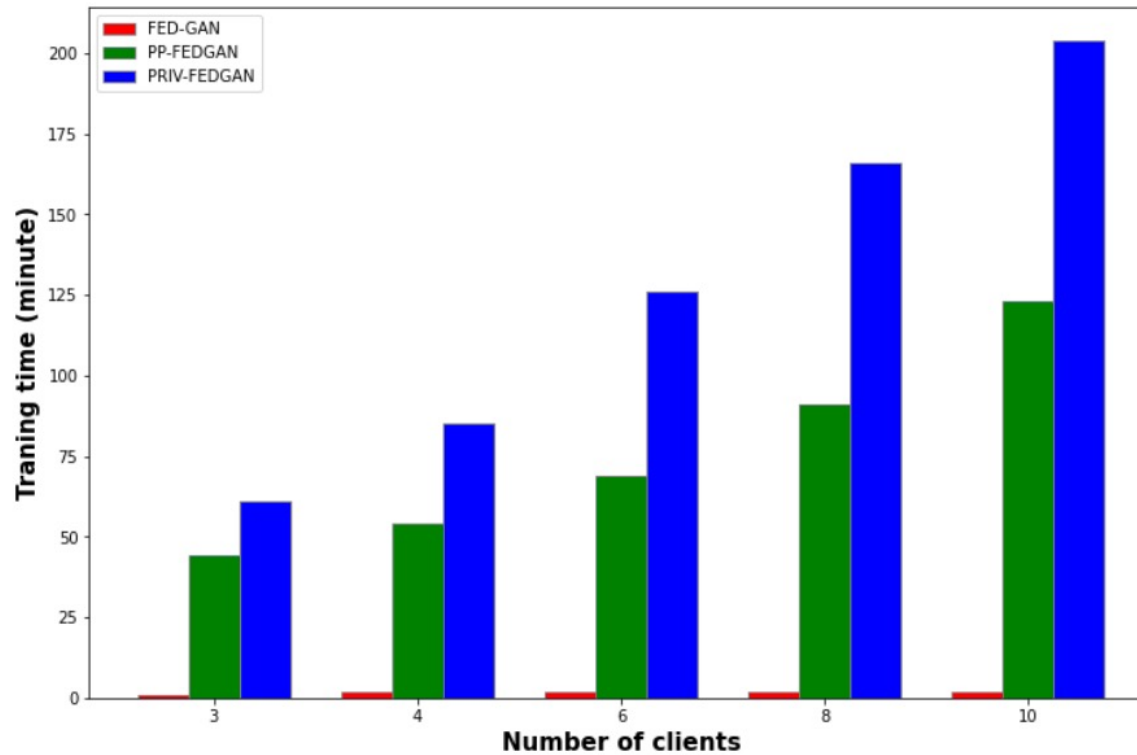
Training overhead

Model	Number of Parameters	Before Enc. (MB)	After Enc. (MB)	Encryption time (s)
Discriminator	109440	0.432	708.34	6.1
Generator	312256	1.2	708.79	10.2

	Each round (Client-side)	Each round (Server-side)	10 rounds training
FedGAN	26	-	261
FedGAN + HE	36	115	1531
FedGAN + DP	78	-	784
PP-FedGAN	89	115	2096



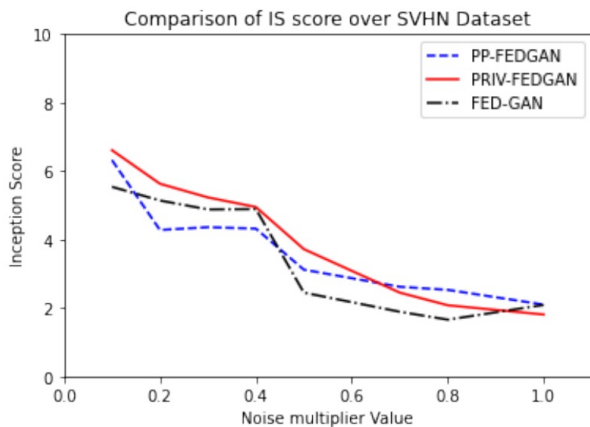
Performance overhead



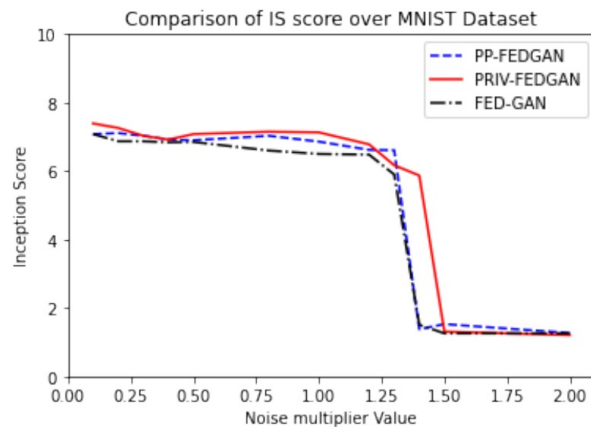


Evaluating the output quality

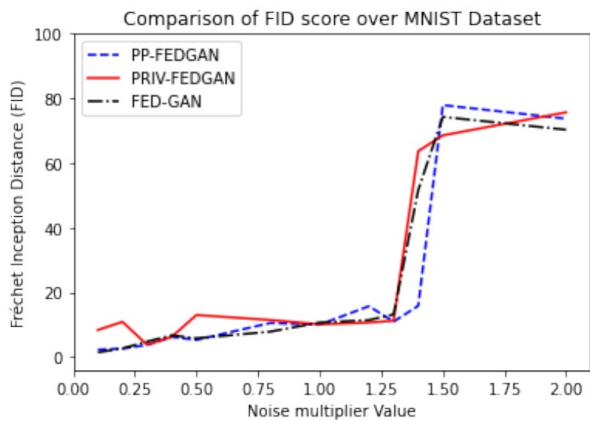
- › The Inception Score (IS)
- › Fréchet Inception Distance (FID)
- › Kernel Inception Distance (KID)



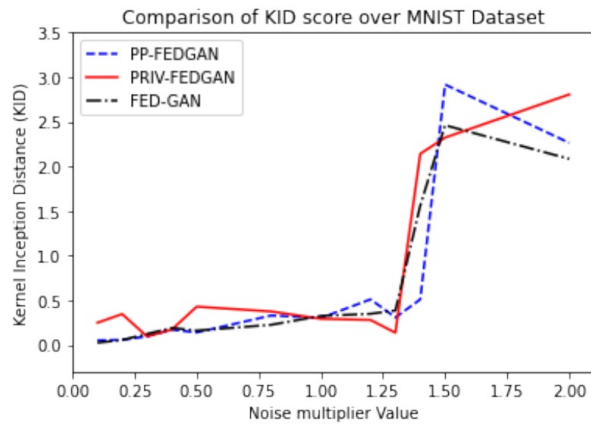
(a)



(b)



(c)



(d)



Conclusion

- › A hybrid privacy-preserving synthetic data generation framework
- › Utilizes homomorphic encryption and differential privacy
- › Guarantees against semi-honest adversaries
- › High quality output
- › 30% better training performance than existing work



References

- › [1] C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor, “Our data, ourselves: Privacy via distributed noise generation,” in Annual International Conference on the Theory and Applications of Cryptographic Techniques. Springer, 2006, pp. 486–503.
- › [2] McMahan, Brendan, et al. "Communication-efficient learning of deep networks from decentralized data." Artificial intelligence and statistics. PMLR, 2017.
- › [3] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep learning with differential privacy,” in Proceedings of the 2016 ACM SIGSAC conference on computer and communications security, 2016, pp. 308–318.
- › [4] Creswell, Antonia, et al. "Generative adversarial networks: An overview." IEEE signal processing magazine 35.1 (2018): 53-65.
- › [5] R. C. Geyer, T. Klein et al., “Differentially private federated learning: A client level perspective,” arXiv preprint arXiv:1712.07557, 2017.
- › [6] Xin, Bangzhou, et al. "Federated synthetic data generation with differential privacy." Neurocomputing 468 (2022): 1-10.
- › [7] Cheon, Jung Hee, et al. "Homomorphic encryption for arithmetic of approximate numbers." Advances in Cryptology–ASIACRYPT 2017: 23rd International Conference on the Theory and Applications of Cryptology and Information Security, Hong Kong, China, December 3-7, 2017, Proceedings, Part I 23. Springer International Publishing, 2017.



a.r.ghavamipour@rug.nl