



An Introduction to Federated Learning

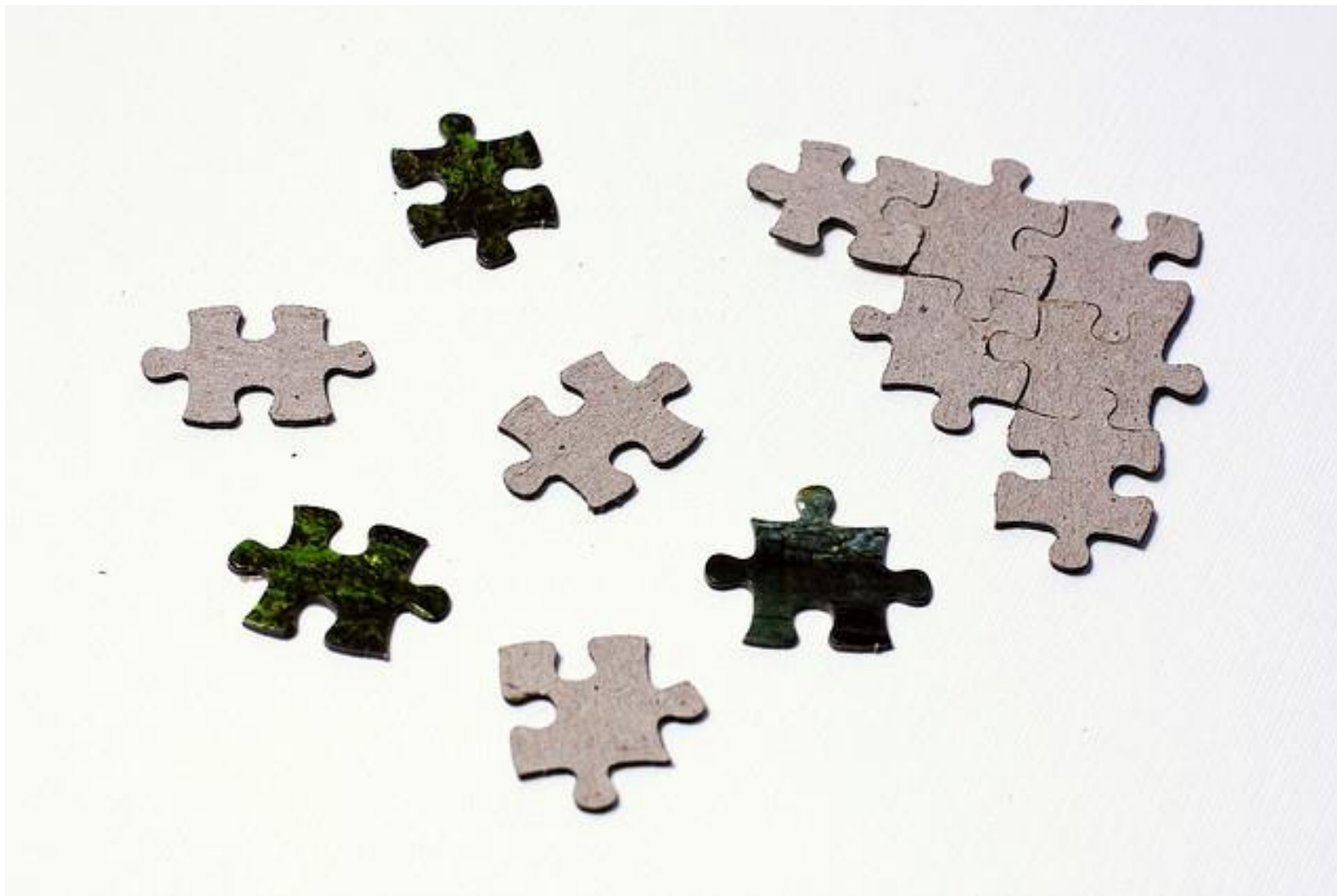
Qiang Yang

WeBank, HKUST



<https://www.fedai.org/>

Challenge for AI : Data Fragmentation, Data Silos



Background

- Increasingly strict laws on data protection:
 - GDPR of EU, 2018
 - CCPA of USA, 2018
 - Cyber Security Law of China, 2017
- Growing concern on user privacy and data security
- Data exist in the form of isolated silos.
- **Federated learning can be a solution!** [McMahan'16, Yang'19]
- Reference:
 - [1] DLA Piper, Data Protection Laws of the World <https://www.dlapiperdataprotection.com/>
 - [2] DLA Piper, Data Protection Laws of the World, Full Handbook

DLA Piper



Image from DLA Piper

Data Sharing Among Parties: Difficult, Impossible, Illegal

- Medical clinical trial data cannot be shared (by R. Stegeman 2018 on Genemetics)
- Our society demands more control on data privacy and security
 - GDPR, Government Regulations
 - Corporate Security and Confidentiality Concerns
 - Data privacy concerns



China's Data Privacy Laws

- Many enacted since 2017
- Requires that Internet businesses must not leak or tamper with the personal information
- When conducting data transactions with third parties, they need to ensure that the proposed contract follow legal data protection obligations.
- More to come...

From Report by KPMG 2017

Highlights and interpretation of the Cybersecurity Law



Highlights of the Cybersecurity Law

Comprising 79 articles in seven chapters, the Cybersecurity Law contains a number of cybersecurity requirements, including safeguards for national cyberspace sovereignty, protection of critical information infrastructure and data and protection of individual privacy. The Law also specifies the cybersecurity obligations for all parties. Enterprises and related organisations should prioritise the following highlights of the Cybersecurity Law:



Personal information protection

The Cybersecurity Law clearly states requirements for the collection, use and protection of personal information.



Critical information infrastructure

The Cybersecurity Law frequently mentions the protection of "critical information infrastructure".



Network operators

"Network operators" are the owners and administrators of networks and network service providers. The Cybersecurity Law clarifies operators' security responsibilities.



Preservation of sensitive information

The Cybersecurity Law requires personal information/important data collected or generated in China to be stored domestically.



Certification of security products

Critical cyber equipment and special cybersecurity products can only be sold or provided after receiving security certifications.



Legal liabilities

Enterprises and organisations that violate the Cybersecurity Law may be fined up to RMB1,000,000.

Facebook finally rolls out privacy tool for your browsing history

By Kaya Yurieff, CNN Business
Updated 1839 GMT (0229 HKT) August 2



Google strengthens Chrome's privacy controls

Frederic Lardinois @frederic / 7

Google today announced that will, in the long run, introduce cookies and enhance its user interface. With this move, Google is not anti-fingerprinting technology happening in the Chrome browser change and adapt their code.

Top Microsoft exec says online privacy has reached 'a crisis point'

By Clare Duffy, CNN Business
Updated 1749 GMT (0149 HKT) October 14, 2019



Challenges to AI: small data and fragmented data, Non-iid, Non-balanced, non-cooperative or malicious, dirty data, incomplete data, outdated data...

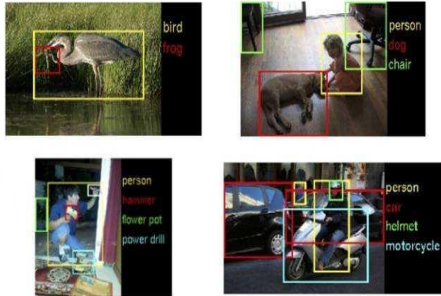


Enterprise A

X1



Data silos



Enterprise B

(X2, Y)

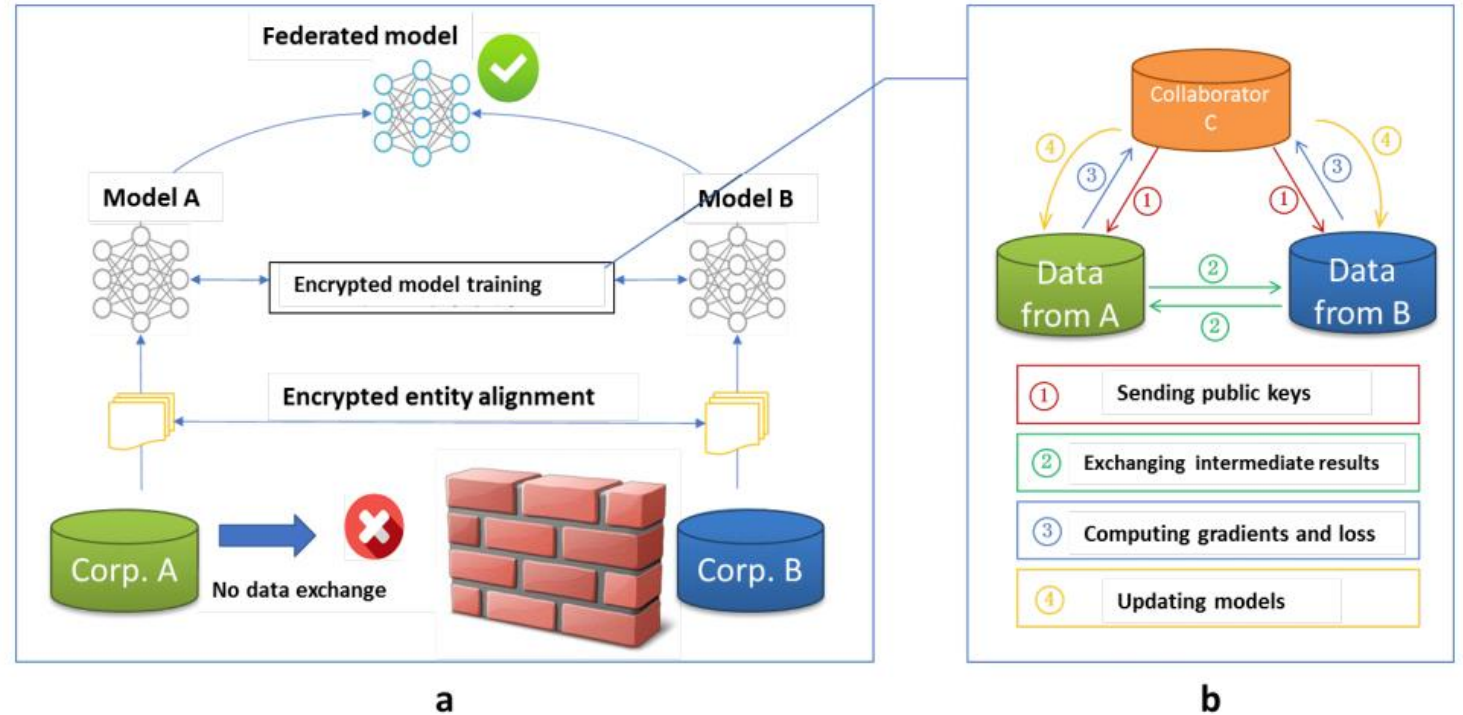
Low Security in Data Sharing
Lack of Labeled Data
Segregated Datasets

Over 80% of the worlds' enterprise information are in data silos!



Federated Learning

- Move models instead of data; Data usable but not visible.
- Multi-party model learning without exchanging data.
- Scenarios:
 - Collaboration among edge devices,
 - Collaboration among organizations,
 - Collaboration among departments within one organization.
- Also known as:
 - Federated Machine Learning
 - Collaborative Machine Learning
 - Federated Deep Learning
 - Federated Optimization
 - Privacy-preserving Machine Learning
 - Geo-distributed Machine Learning
 - Geo-distributed Deep Learning
 - Multi-party Learning

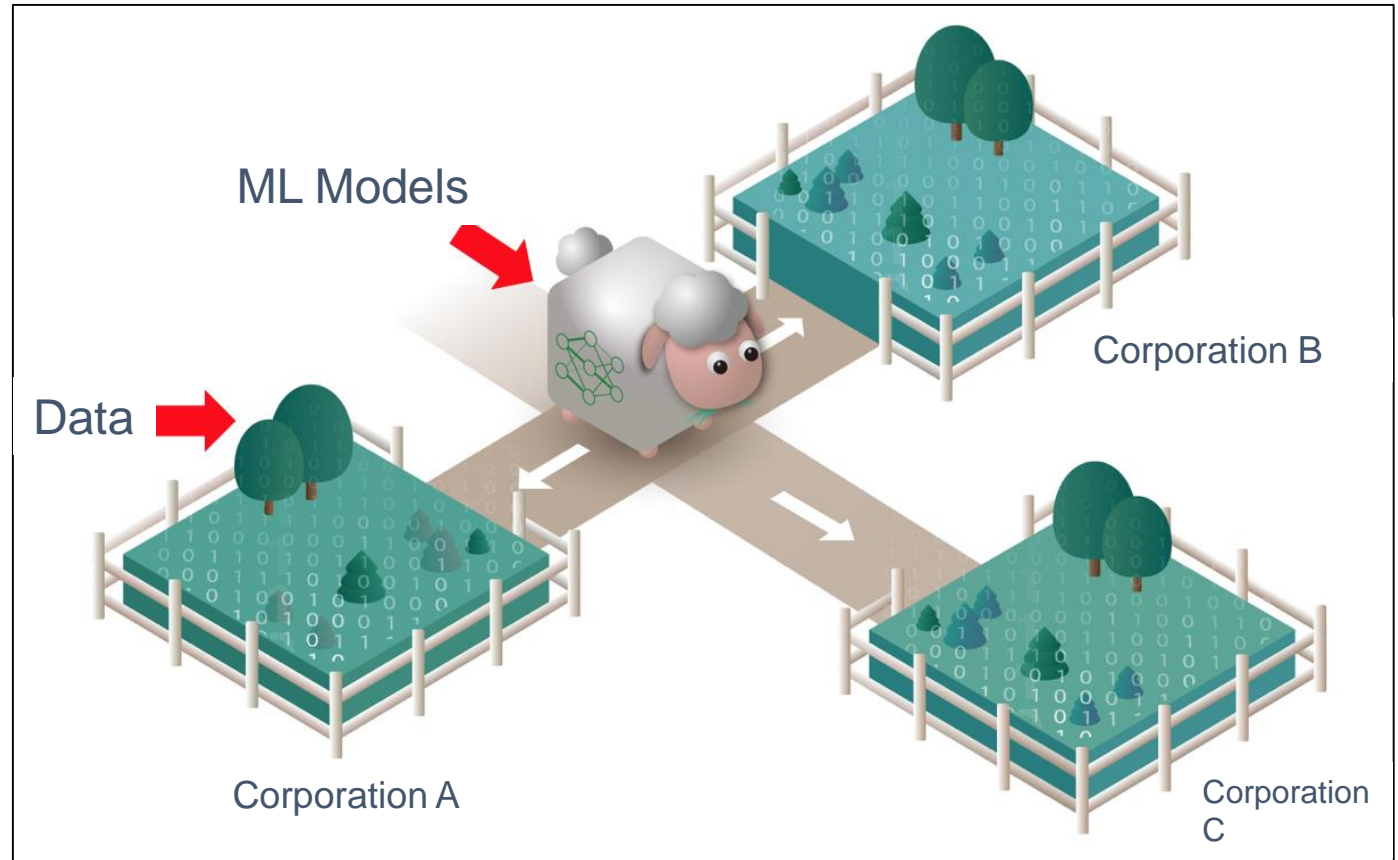


References:

- [Kairouz'19] Peter Kairouz, and H. Brendan McMahan, et. al., "Advances and Open Problems in Federated Learning," Dec. 2019. Available: <https://arxiv.org/abs/1912.04977>
- [Yang'19] Qiang Yang, et al., "Federated Machine Learning: Concept and Applications," Feb. 2019.
- [McMahan'16] H. Brendan McMahan, et al., "Federated Learning of Deep Networks using Model Averaging," Feb. 2016.
- [Konecny'15] Jakub Konečný, Brendan McMahan, and Daniel Ramage, "Federated Optimization: Distributed Optimization Beyond the Datacenter," Nov. 2015. Available: <https://arxiv.org/abs/1511.0357>

Definition of Federated Learning-2

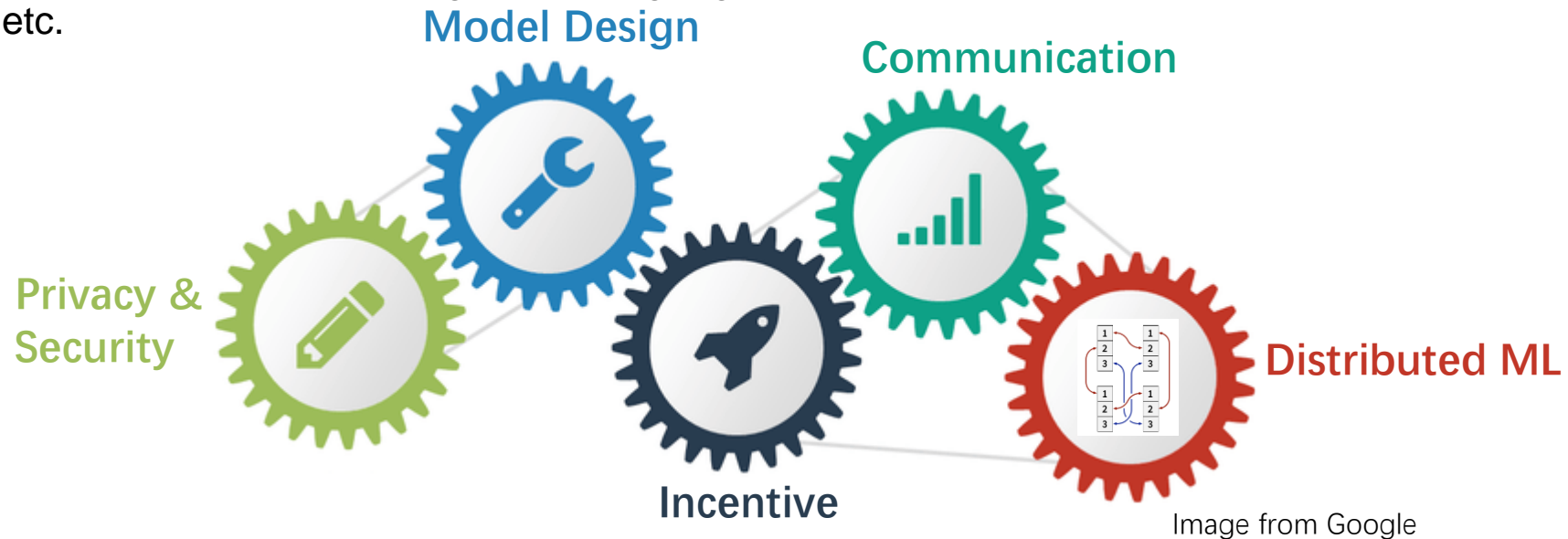
- Interpretation of Federated Learning:
 - Models --- Sheep
 - Data --- Grass
- Originally, one need to purchase grass from different sources to feed sheep --- Companies gather lots of data to train models, where many challenges exist, such as user privacy, data security and regulations.
- Federated Learning provides an alternative: sheep are led to different farms and can thus eat grass from all places without having to move the grass. --- Federated learning models gather knowledge from various sources of data without having to observe them.



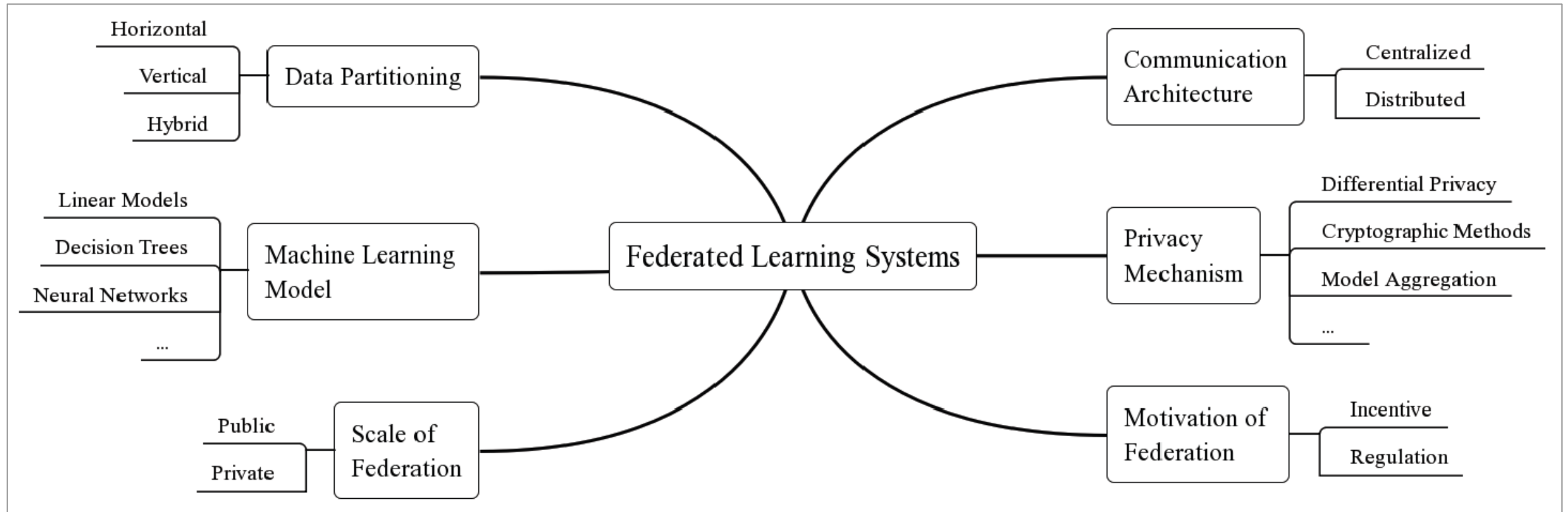
Key Components in Federated Learning

<<Book: Federated Learning>>

- Model design and hyperparameter tuning, e.g. number of layers, CNN or RNN, etc.
- Distributed learning algorithm (Chapter 3), e.g. client selection, tackling non-IID (or even contradictory) training data, system-algorithm co-design, etc.
- Communication optimization, e.g. alleviating the influence of network delay, model/gradient compression, etc.
- Security and privacy (Chapter 2), e.g. Homomorphic Encryption (HE), Differential Privacy (DP), Secure Multi-party Computation (MPC), etc.
- Incentive mechanism (Chapter 7), e.g. motivating organizations from different industries, adequate revenue allocation, etc.



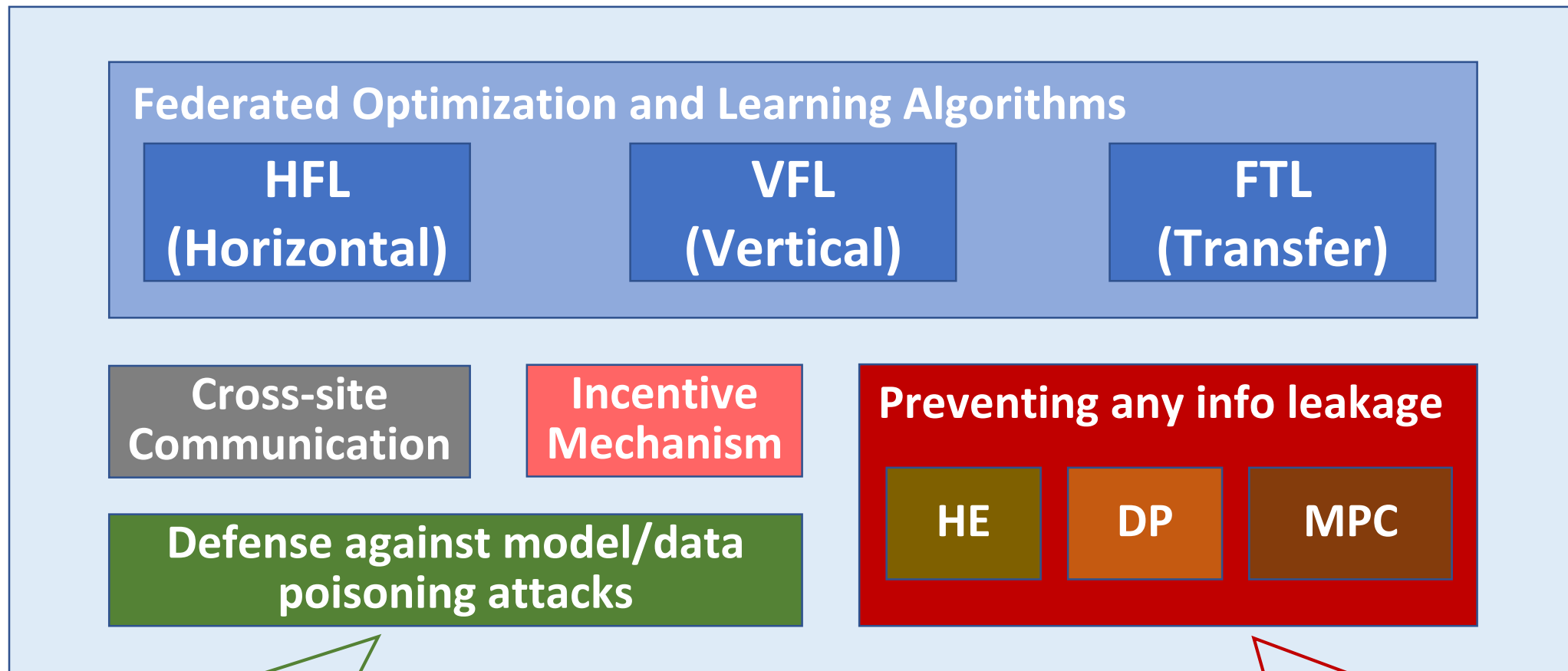
Federated Learning Systems: Overview



Taxonomy of federated learning systems (FLSs)

[Li'19] Qinbin Li, Zeyi Wen, et al., "Federated Learning Systems: Vision, Hype and Reality for Data Privacy and Protection," Oct. 2019. <https://arxiv.org/abs/1907.09693>

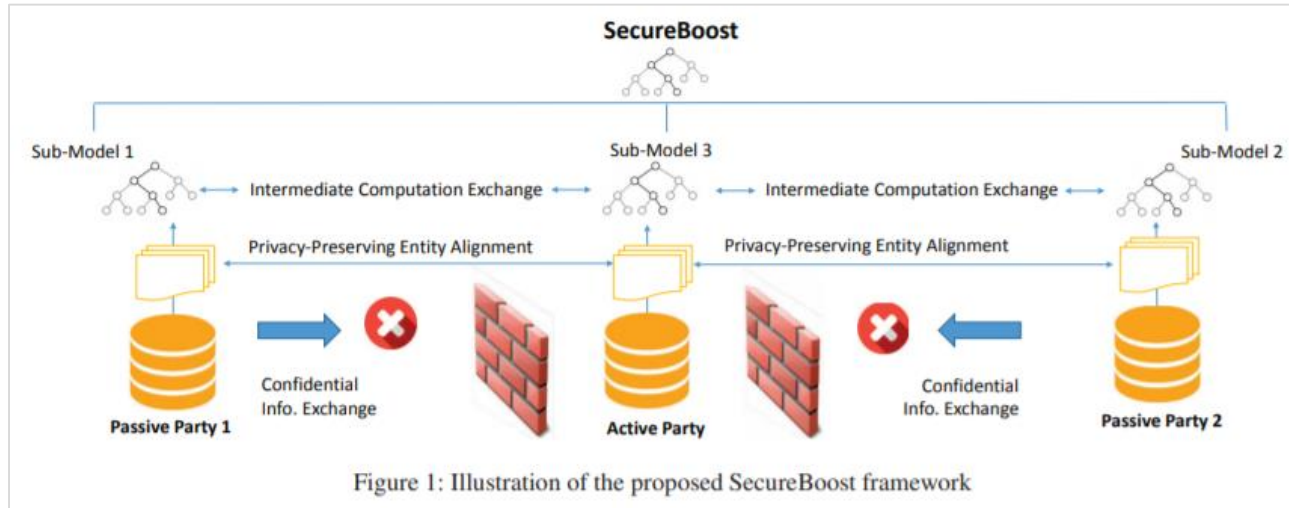
Federated Learning System: Overview



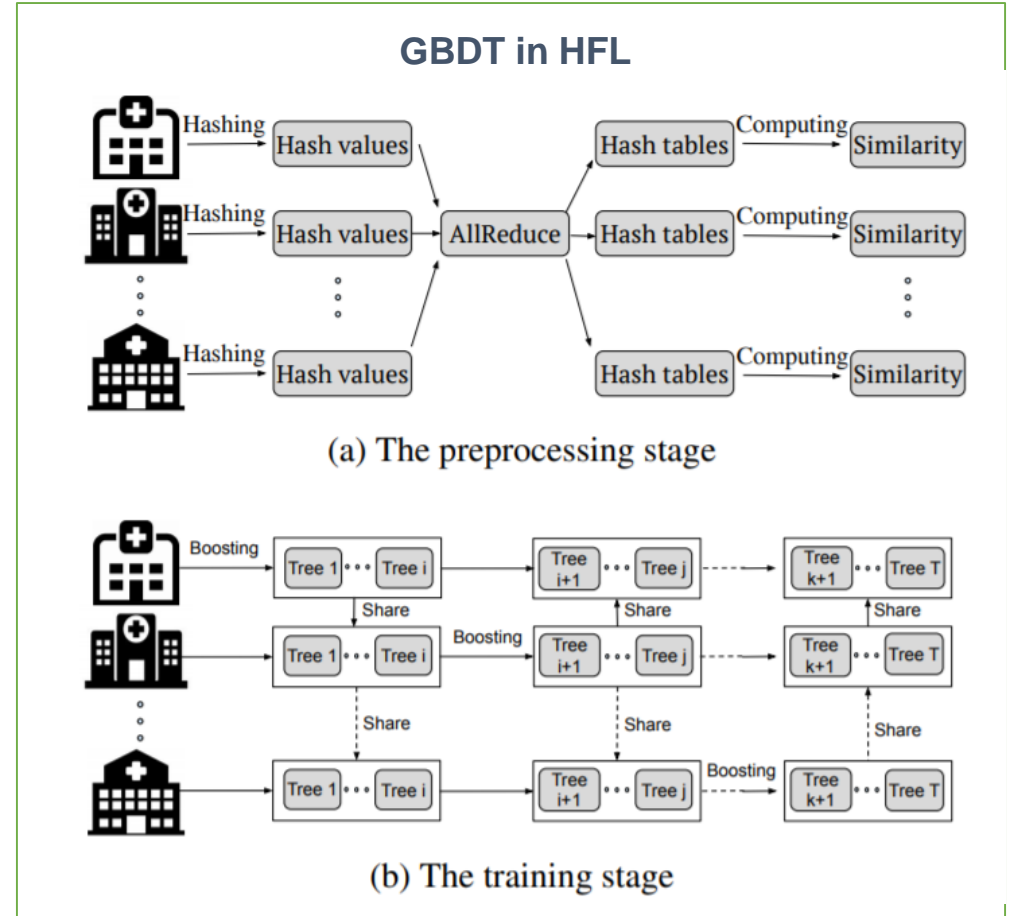
- FL has built-in mechanism for robustness, such as defending model/data poisoning attacks.
- FL is more than what “MPC+ML” is about.

- The built-in MPC block in FL can prevent any info leakage (either model or data leakage).
- FL can do whatever “MPC+ML” can do.

Secureboost in VFL

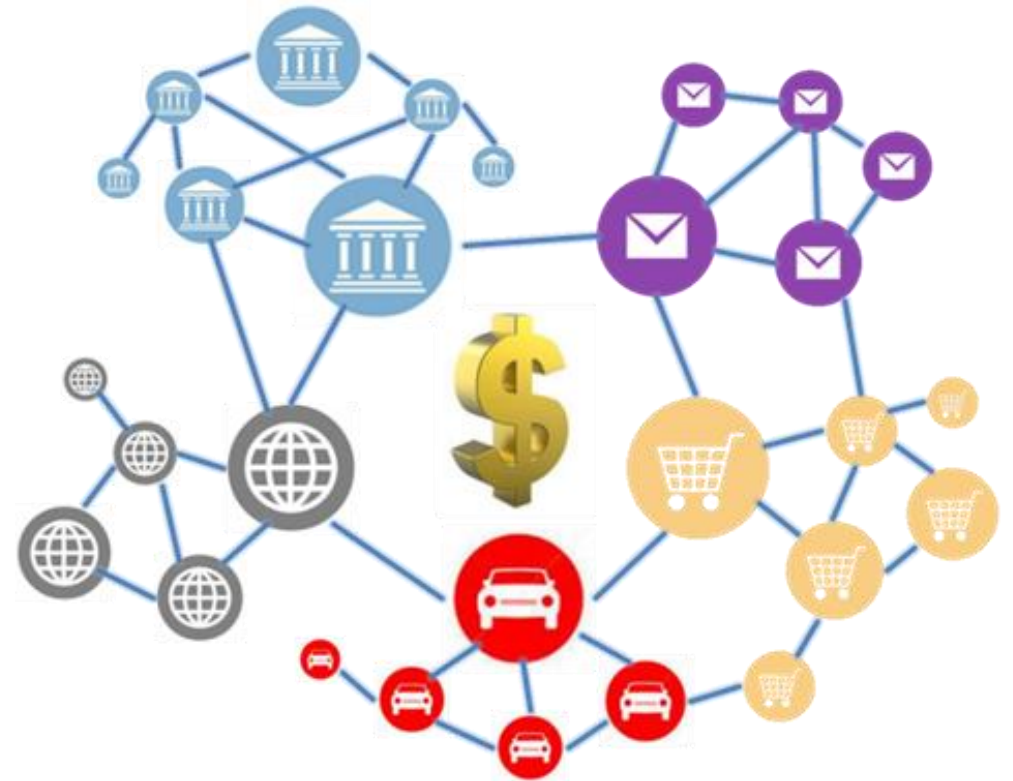
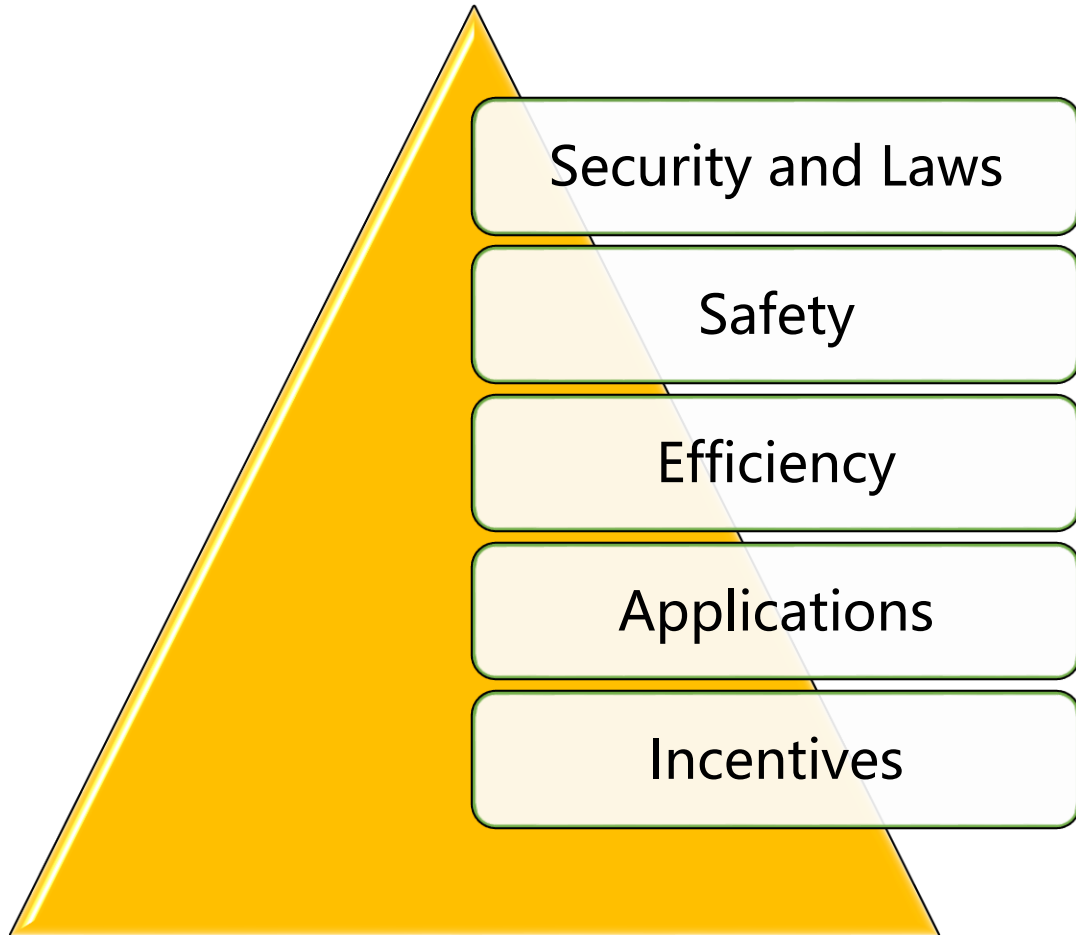


[Kewei Cheng](#), [Tao Fan](#), [Yilun Jin](#), [Yang Liu](#), [Tianjian Chen](#), [Qiang Yang](#),
SecureBoost: A Lossless Federated Learning Framework, IEEE Intelligent Systems 2020

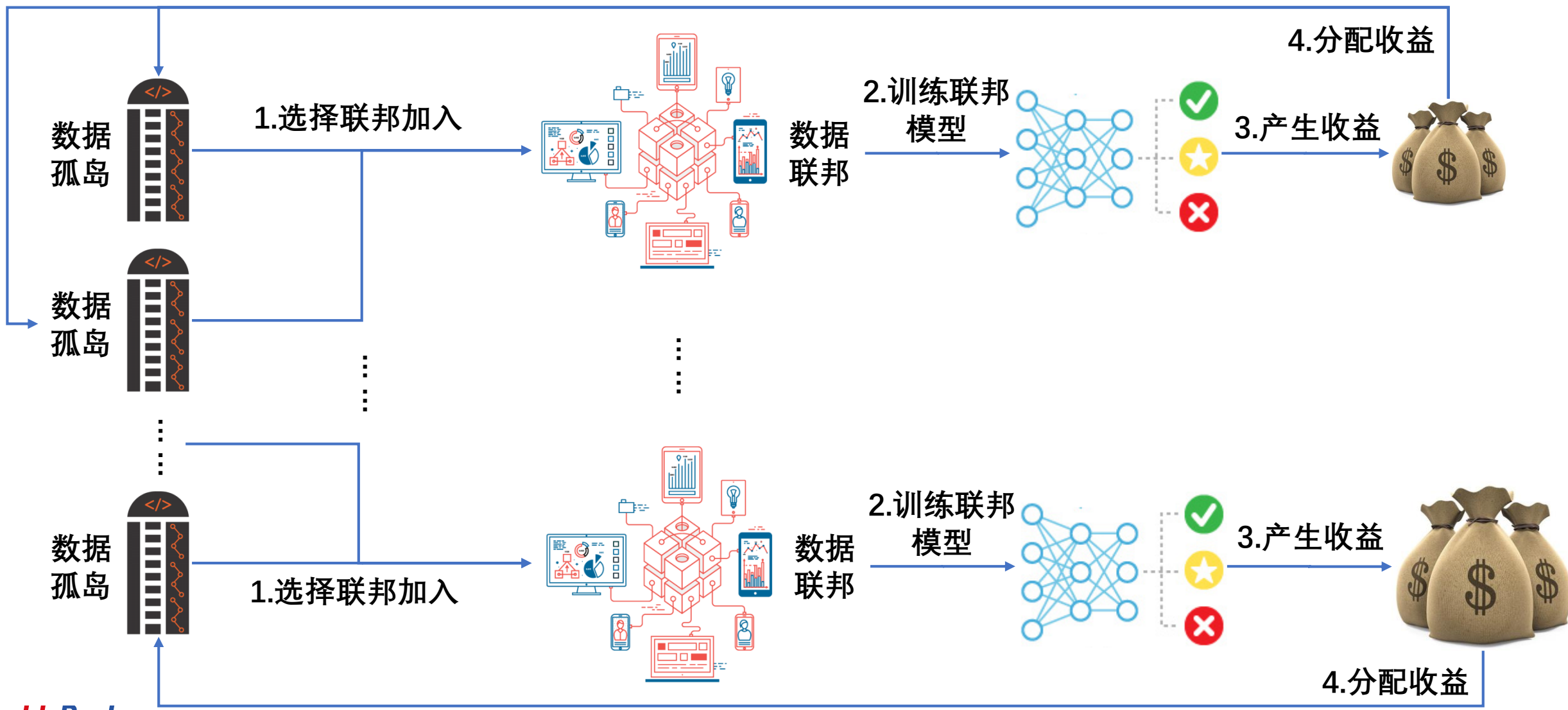


[Qinbin Li](#), [Zeyi Wen](#), [Bingsheng He](#), Practical Federated Gradient Boosting Decision Trees, AAAI, 2019

Federated Learning: research areas



Federated Learning Market Games



Federated Learning Applications

- Federated Learning + Other ML Algorithms (Chapter 8)
 - Federated learning + Computer vision (FL+CV)
 - Federated learning + Natural language processing (FL+NLP), including automatic speech recognition (ASR)
 - Federated learning + Recommender system (FL+RS)
- Federated Learning + Industry/Society, (Chapter 10 and cases from WeBank: <https://www.fedai.org/cases/>)
 - Federated learning + Finance, FinTech
 - Federated learning + Insurance, InsurTech
 - Federated learning + Healthcare
 - Federated learning + Education
 - Federated learning + AIoT
 - Federated learning + Smart City
 - Federated learning + Edge computing
 - Federated learning + 5G/6G
- Reference:
 - GitHub, innovation-cat/Awesome-Federated-Machine-Learning, Available: <https://github.com/innovation-cat/Awesome-Federated-Machine-Learning>
 - <https://zhuanlan.zhihu.com/p/87777798>

Federated Learning Datasets

- WeBank FedVision - Street Dataset, Available: <https://dataset.fedai.org/#/>
 - A real-world **object detection** dataset that annotates images captured by a set of street cameras based on object present in them, including 7 classes. In this dataset, each or every few cameras serve as a device.
- Carnegie Mellon University, LEAF: A benchmarking framework for federated learning, Available: <https://leaf.cmu.edu/>, <https://github.com/TalwalkarLab/leaf>
 - Federated Extended MNIST (FEMNIST), 62 classes, Image Classification
 - Twitter, Sentiment140, Sentiment Analysis, federated
 - Shakespeare, Next-Character Prediction, federated
 - Celeba, Image Classification (Smiling vs. Not smiling), federated
 - Synthetic Dataset, Classification, federated
- University of Southern California: FedML (Chaoyang He et al.)

Federated Learning Open-source Platforms

- **WeBank FATE**, supports TensorFlow and PyTorch, <https://github.com/FederatedAI/FATE>
- Google TensorFlow Federated (TFF) , <https://github.com/tensorflow/federated>
- Google TensorFlow-Encrypted, <https://github.com/tf-encrypted/tf-encrypted>
- PyTorch, torch.distributed, https://pytorch.org/tutorials/intermediate/dist_tuto.html
- **Uber Horovod**, supports Keras, TensorFlow, PyTorch, and MXNet, <https://github.com/horovod/horovod>
- coMindOrg, supports TensorFlow, <https://github.com/coMindOrg/federated-averaging-tutorials>
- OpenMined PySyft, supports PyTorch, <https://github.com/OpenMined/PySyft>
- MesaTEE by Baidu, <https://mesatee.org/> <https://mp.weixin.qq.com/s/1SXW1N7BaVnyFXFZ-f24TA>
- PaddlePaddle/PaddleFL by Baidu, <https://github.com/PaddlePaddle/PaddleFL>

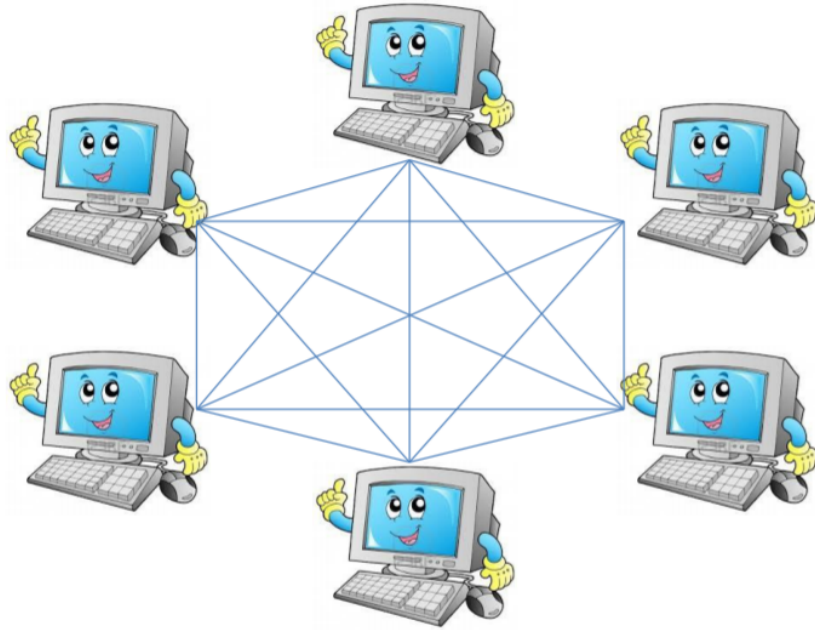
Privacy-Preserving Technologies

- Secure Multi-party Computation (MPC)
- Homomorphic Encryption (HE)
- Yao's Garbled Circuit
- Secret sharing
- Differential Privacy (DP)

.....



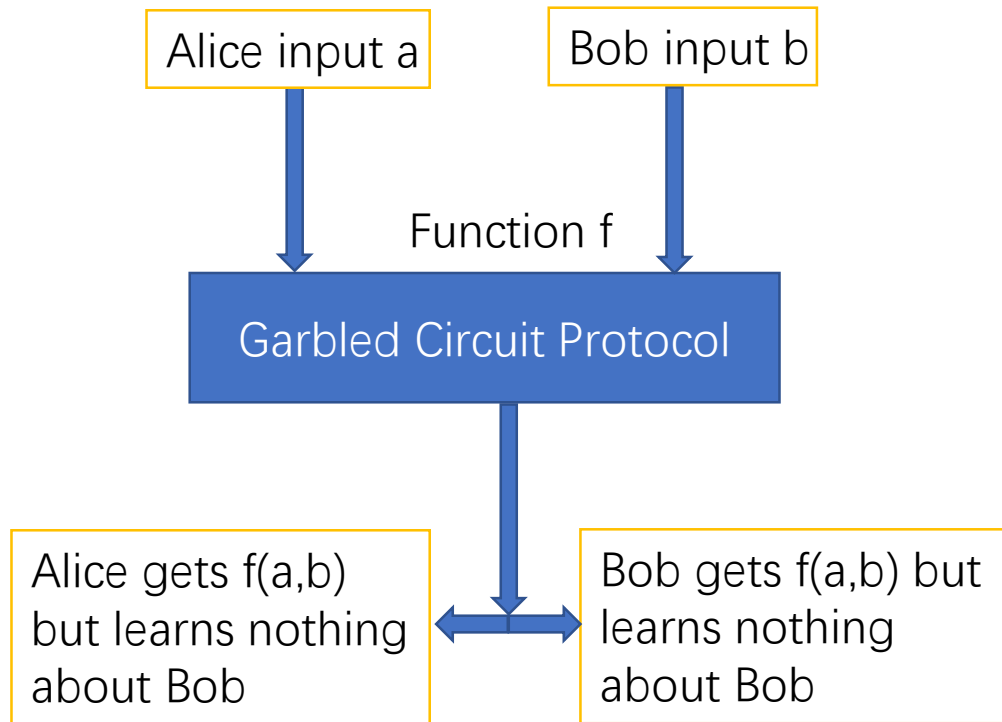
Secure Multi-Party Computation (MPC)



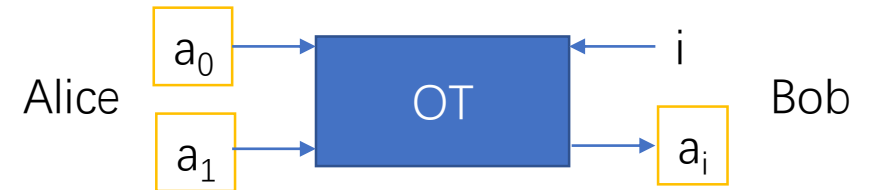
Ran Cohen, Tel Aviv University, Secure Multiparty Computation: Introduction

- Provides security proof in a well-defined simulation framework
- Guarantees complete zero knowledge
- Requires participants' data to be secretly-shared among non-colluding servers
- Drawbacks:
 - Expensive communication,
 - Though it is possible to build a security model with MPC under lower security requirement in exchange for efficiency

Yao's Garbled Circuit Protocol (Andrew Yao, 1986)



- Oblivious Transfer

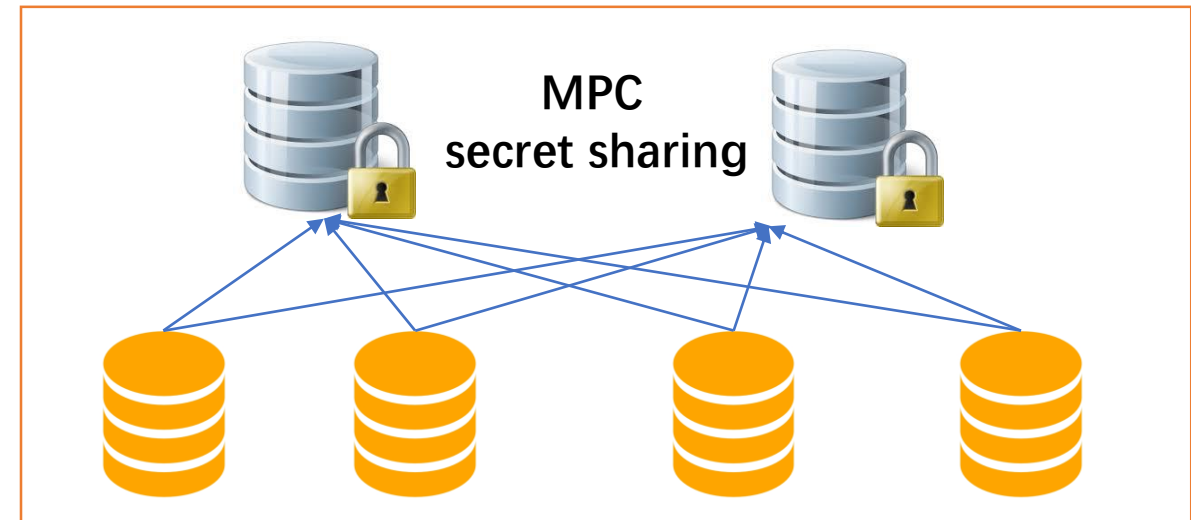


Steps

- Alice builds a garbled circuits;
- Alice sends her input keys;
- Alice and Bob do Oblivious Transfer;
- Bob gets the output and sends back to Alice;
- Alice and Bob learns nothing about the other value.

SecureML: Privacy-preserving machine learning for linear regression, logistic regression and neural network training

- Combines secret sharing, garbled circuits and oblivious transfer
- Learns via two un-trusted, but non-colluding servers
- Computationally expensive



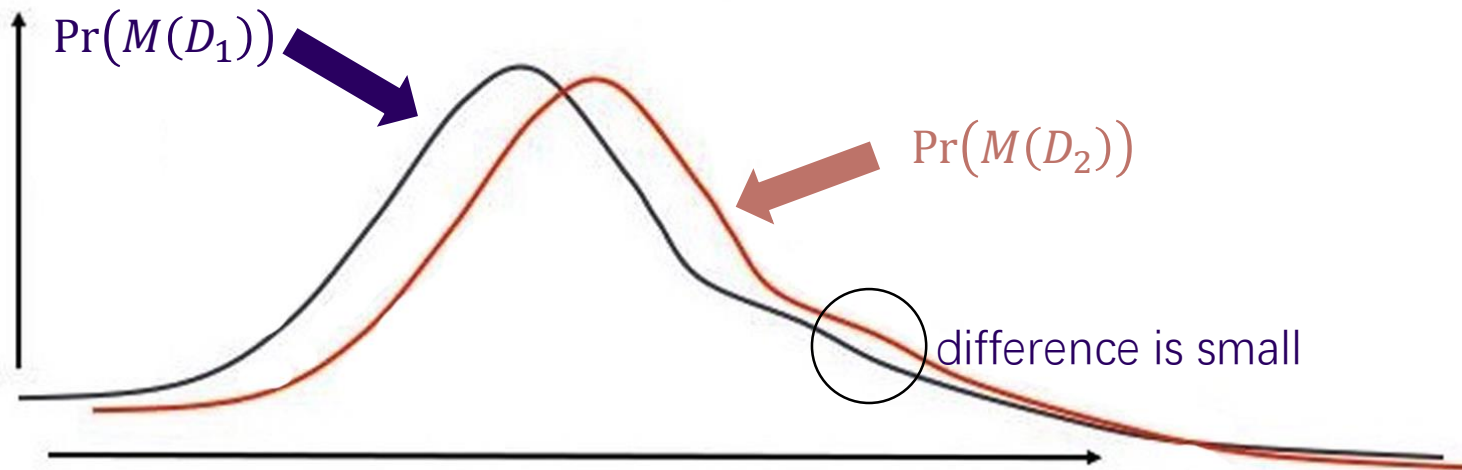
Mohassel, P., & Zhang, Y. (2017, May). SecureML: A system for scalable privacy-preserving machine learning. In *2017 38th IEEE Symposium on Security and Privacy (SP)* (pp. 19-38). IEEE.

Differential Privacy

Definition: Differential Privacy (DP) [Dwork 2008]

A randomized mechanism M is ϵ -differentially private, if for all output t of M , and for all databases D_1 and D_2 which differ by at most one element, we have

$$\Pr(M(D_1) = t) = e^\epsilon \Pr(M(D_2) = t).$$



Intuition: changes in the distribution are too small to be perceived with variations on a single element.

Cynthia Dwork, 2008. Differential privacy: a survey of results. Theory and Applications of Models of Computation.

Homomorphic Encryption

- Full Homomorphic Encryption and Partial Homomorphic Encryption.
- **Paillier** partially homomorphic encryption

$$\begin{aligned} \textit{Addition} : \quad & [[u]] + [[v]] = [[u+v]] \\ \textit{Scalar multiplication} : \quad & n[[u]] = [[nu]] \end{aligned}$$

- For public key $pk = n$, the encoded form of $m \in \{0, \dots, n - 1\}$ is

$$\text{Encode}(m) = r^n (1 + n)^m \bmod n^2$$

r is randomly selected from $\{0, \dots, n - 1\}$.

- For float $q = (s, e)$, encrypt $[[q]] = ([[s]], e)$, here $q = s\beta^e$ is base- β exponential representation.

Rivest, R. L.; Adleman, L.; and Dertouzos, M. L. 1978. On data banks and privacy homomorphisms. Foundations of Secure Computation, Academia Press 169–179.

Applying HE to Machine Learning

Polynomial approximation for logarithm function

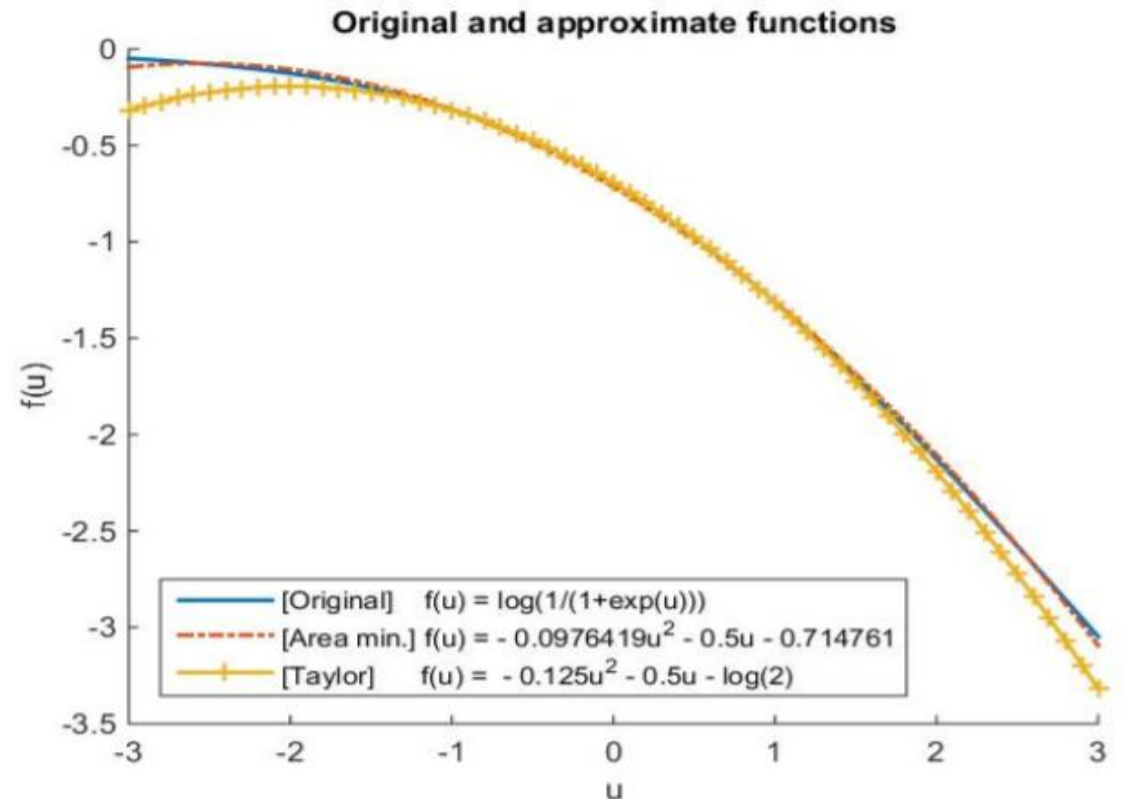
$$\log\left(\frac{1}{1+\exp(u)}\right) \approx \sum_{j=0}^k a_j u^j$$

Encrypted computation for each term in the polynomial function

$$loss = \log 2 - \frac{1}{2} y w^T x + \frac{1}{8} (w^T x)^2$$

$$[[loss]] = [[\log 2]] + \left(-\frac{1}{2}\right) * [[y w^T x]] + \frac{1}{8} [[(w^T x)^2]]$$

- Kim, M.; Song, Y.; Wang, S.; Xia, Y.; and Jiang, X. 2018. Secure logistic regression based on homomorphic encryption: Design and evaluation. JMIR Med Inform 6(2)
- Y. Aono, T. Hayashi, T. P. Le, L. Wang, Scalable and secure logistic regression via homomorphic encryption, CODASPY16



Is the Gradient Info Safe to Share?

Protect gradients with Homomorphic Encryption

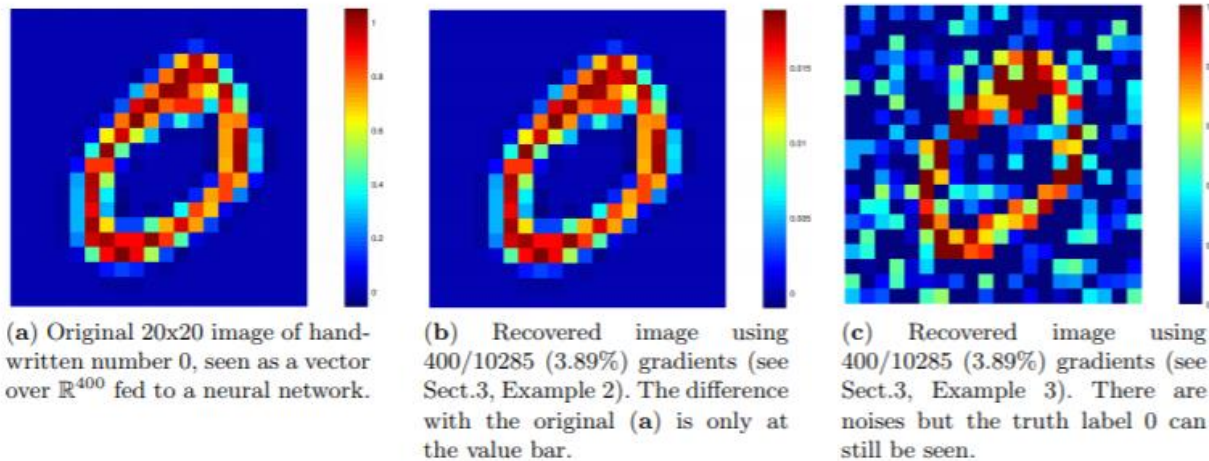
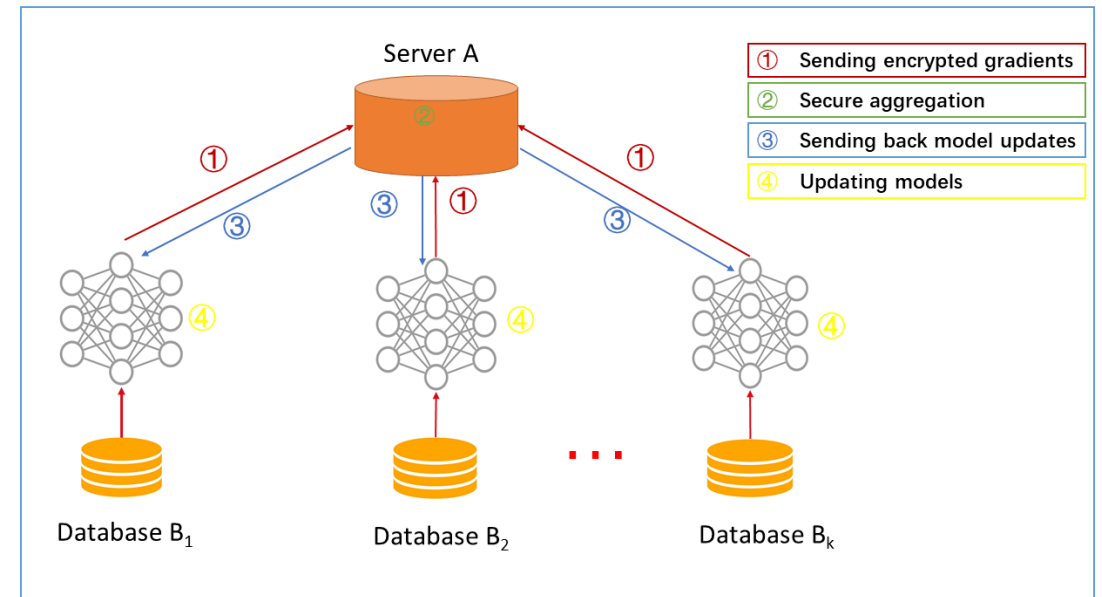


Fig. 3. Original data (a) vs. leakage information (b), (c) from a small part of gradients in a neural network.



Algorithm ensures that no information is leaked to the semi-honest server, provided that the underlying additively homomorphic encryption scheme is secure*.

* Q. Yang, Y. Liu, T. Chen, Y. Tong, Federated machine learning: concepts and applications, ACM TIST, 2018

Le Trieu Phong, et al. 2018. Privacy-Preserving Deep Learning via Additively Homomorphic Encryption. IEEE Trans. Information Forensics and Security, 13, 5 (2018),1333–1345

Horizontal Federated Learning, HFL

- Parties own data with overlapping features, i.e. aligned feature space; yet the training samples are different.
 - Also known as “cross-sample federated learning”, “feature-aligned federated learning”.
 - The feature space is identical.
 - HFL expands the number of training samples, with the feature dimensionality unchanged.

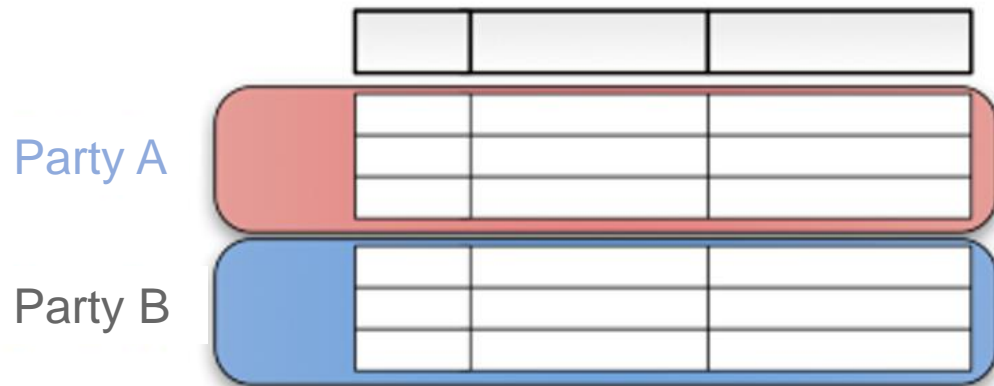
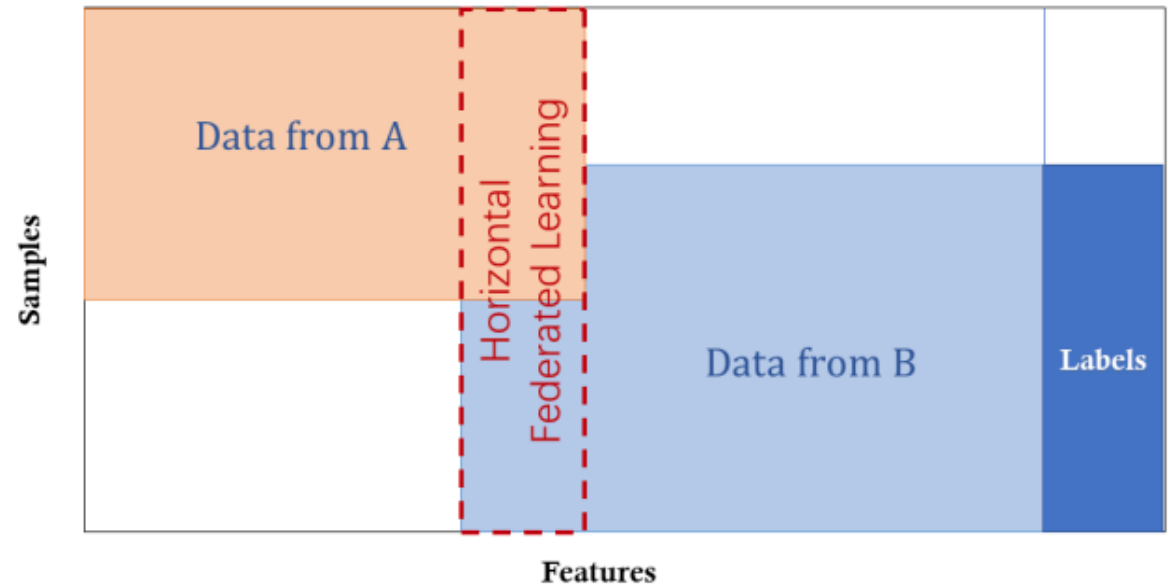


Image from Google

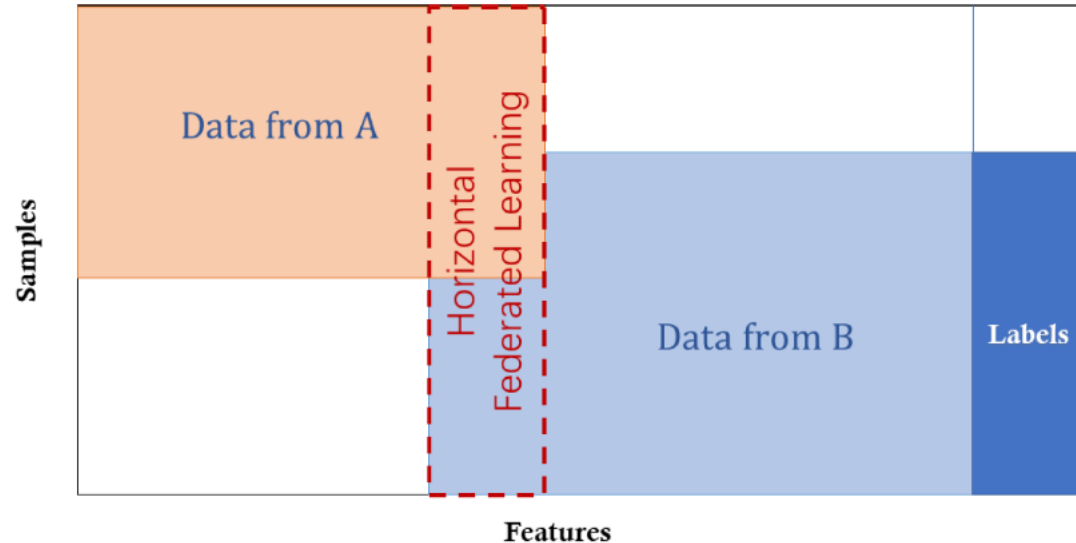
Horizontally partitioned data: data frames are partitioned horizontally into rows, each of which having the same features.



References:

- [Kairouz'19] Peter Kairouz, and H. Brendan McMahan, et. al., "Advances and Open Problems in Federated Learning," Dec. 2019. Available: <https://arxiv.org/abs/1912.04977>
- [Yang'19] Qiang Yang, et al., *Federated machine learning: Concept and Applications*, WeBank, 2019.
- [Google'19] Google Federated Learning Comic, <https://federated.withgoogle.com/>

Horizontal Federated Learning: Divide by Users/Samples



(a) Horizontal Federated Learning

FEDERATED LEARNING FOR MOBILE KEYBOARD PREDICTION, Andrew Hard, et al., Google, 2018

Step 1: Participants compute training gradients locally

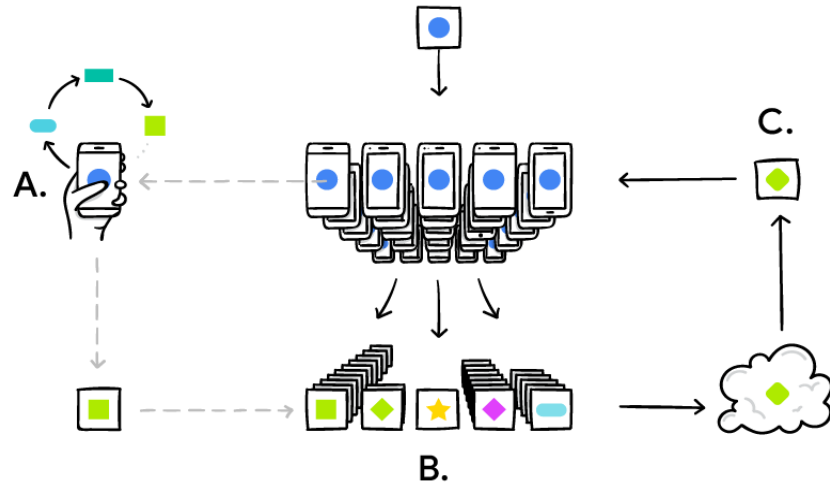
- mask gradients with encryption, differential privacy, or secret sharing techniques
- all participants send their masked results to server

Step 2: The server performs secure aggregation without learning information about any participant

Step 3: The server sends back the aggregated results to participants

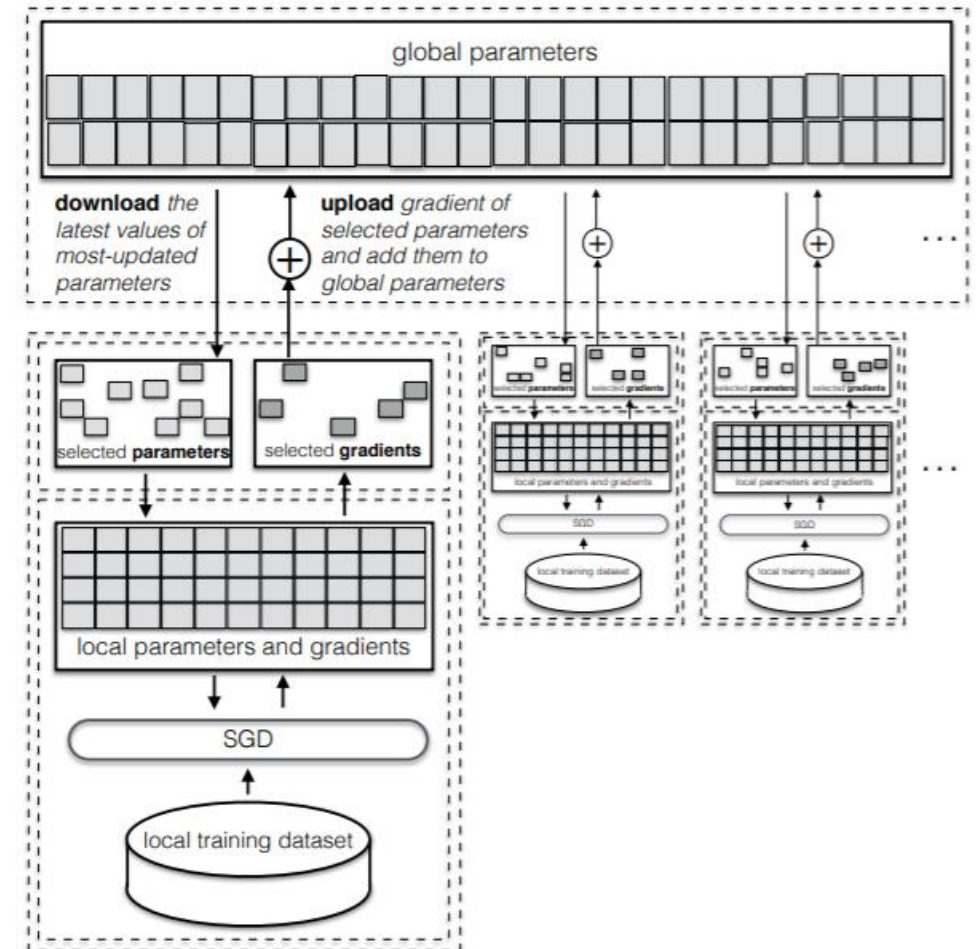
Step 4: Participants update their respective model with the decrypted gradients

Horizontal Federated Learning



H. Brendan McMahan et al, *Communication-Efficient Learning of Deep Networks from Decentralized Data*, Google, 2017

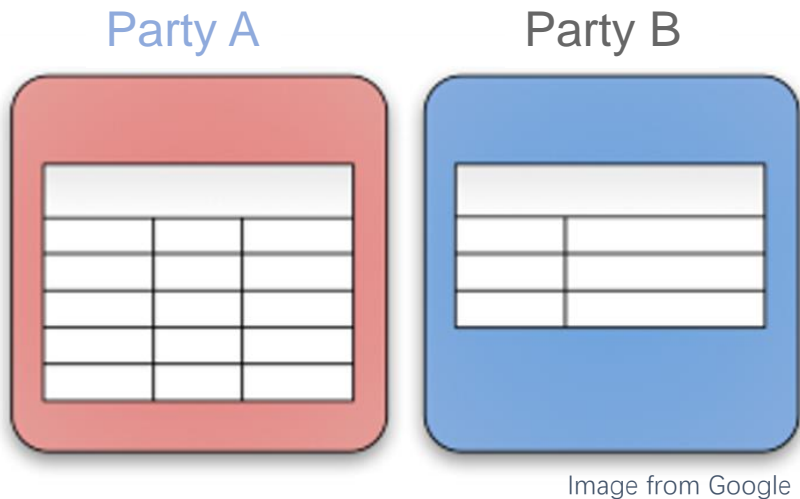
- Multiple clients, one server
- Data is horizontally split across devices, homogeneous features
- Local training
- Selective clients



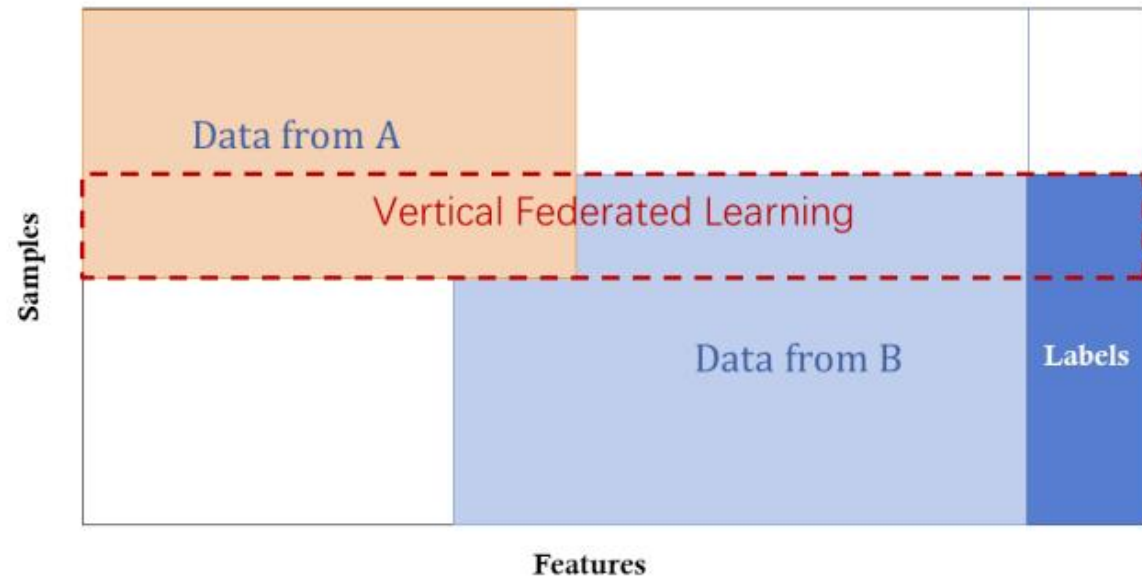
Reza Shokri and Vitaly Shmatikov. 2015. *Privacy-Preserving Deep Learning*. In Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security (CCS '15). ACM, New York

Vertical Federated Learning, VFL

- Parties hold data with identical data ID (i.e. training samples), but with different features.
 - A.k.a “Cross-feature federated learning”, “sample-aligned federated learning”. Suitable for federated learning across industries.
 - Before training, we take the intersection of data IDs held by different parties.
 - VFL increases data dimensionality at the cost of sample size (due to intersection of IDs).



Vertically partitioned data: partition data frames into columns, with each column holding the same feature.



[Yang'19] Qiang Yang, et al., *Federated machine learning: Concept and Applications*, WeBank, 2019.

Vertical Federated Learning

Objective:

- Party (A) and Party (B) co-build a FML model

Assumptions:

- Only one party has label Y
- Neither party wants to expose their X or Y

Challenges:

- Parties with only X cannot build models
- Parties cannot exchange raw data by law

Expectations:

- Data privacy for both parties
- model is LOSSLESS



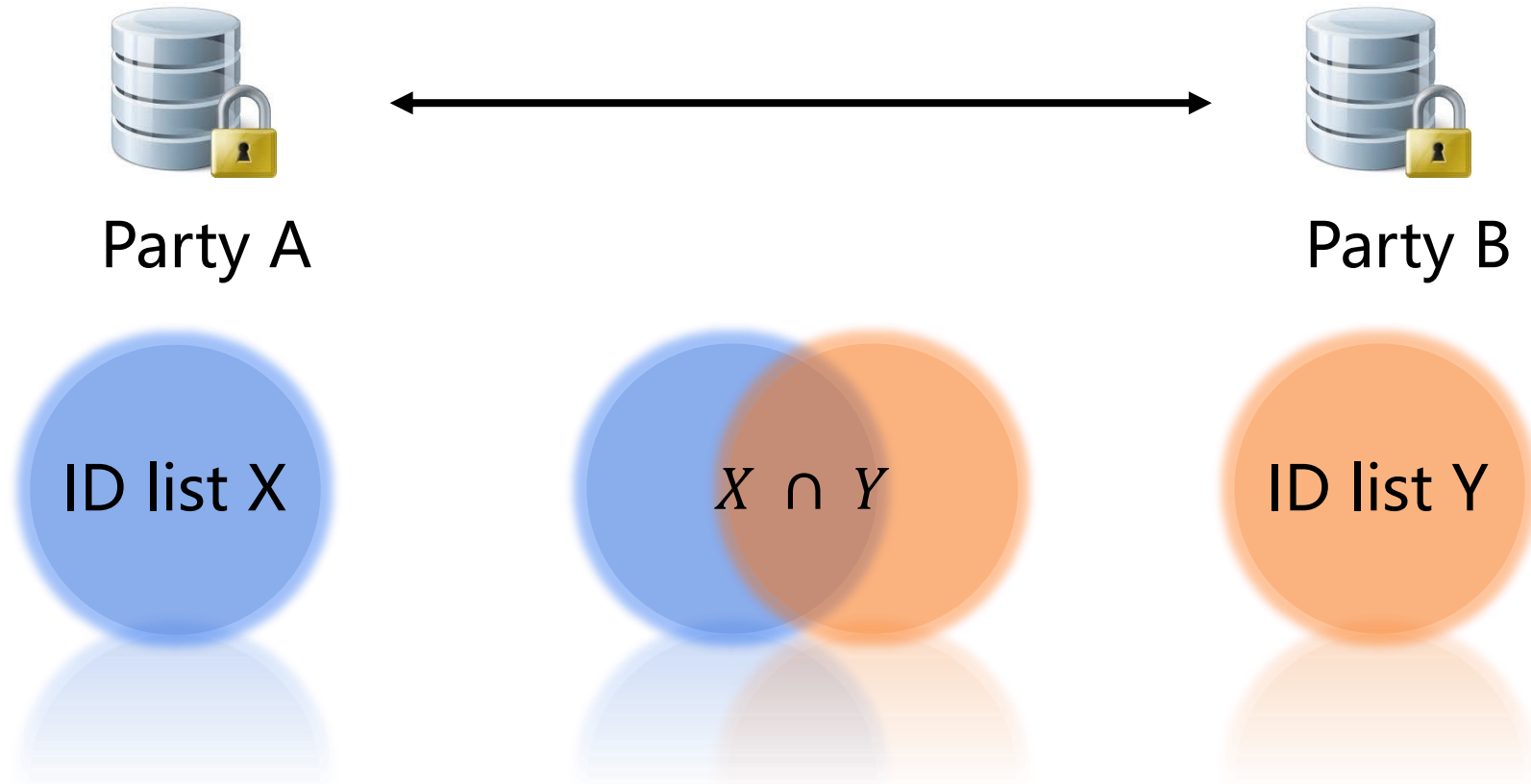
ID	X1	X2	X3	ID	X4	X5	Y
U1	9	80	600	U1	6000	600	No
U2	4	50	550	U2	5500	500	Yes
U3	2	35	520	U3	7200	500	Yes
U4	10	100	600	U4	6000	600	No
U5	5	75	600	U8	6000	600	No
U6	5	75	520	U9	4520	500	Yes
U7	8	80	600	U10	6000	600	No

Retail A Data

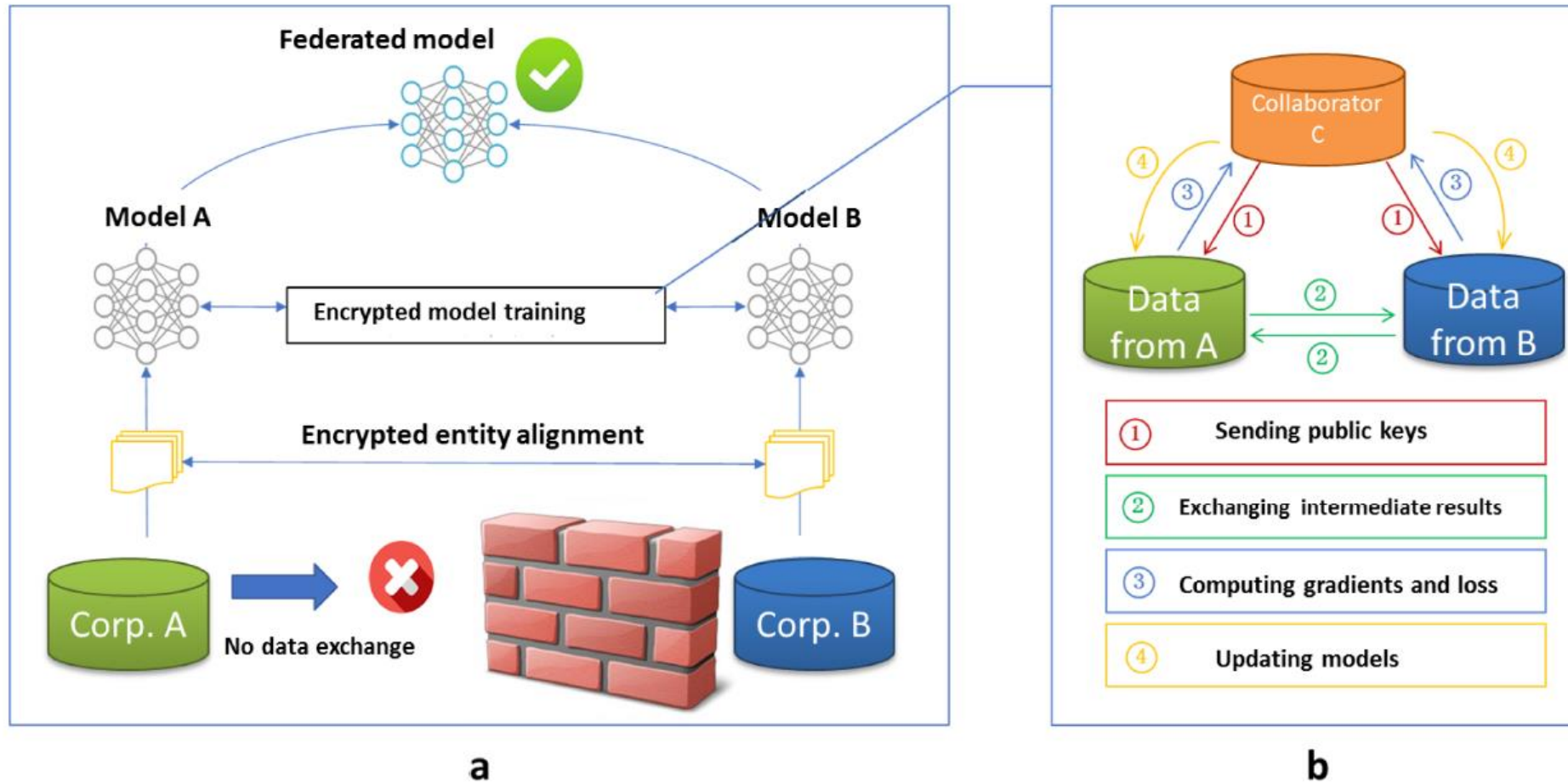
Bank B Data

Privacy-Preserving Entity Match

- Party A and B learns their overlapping IDs but nothing else



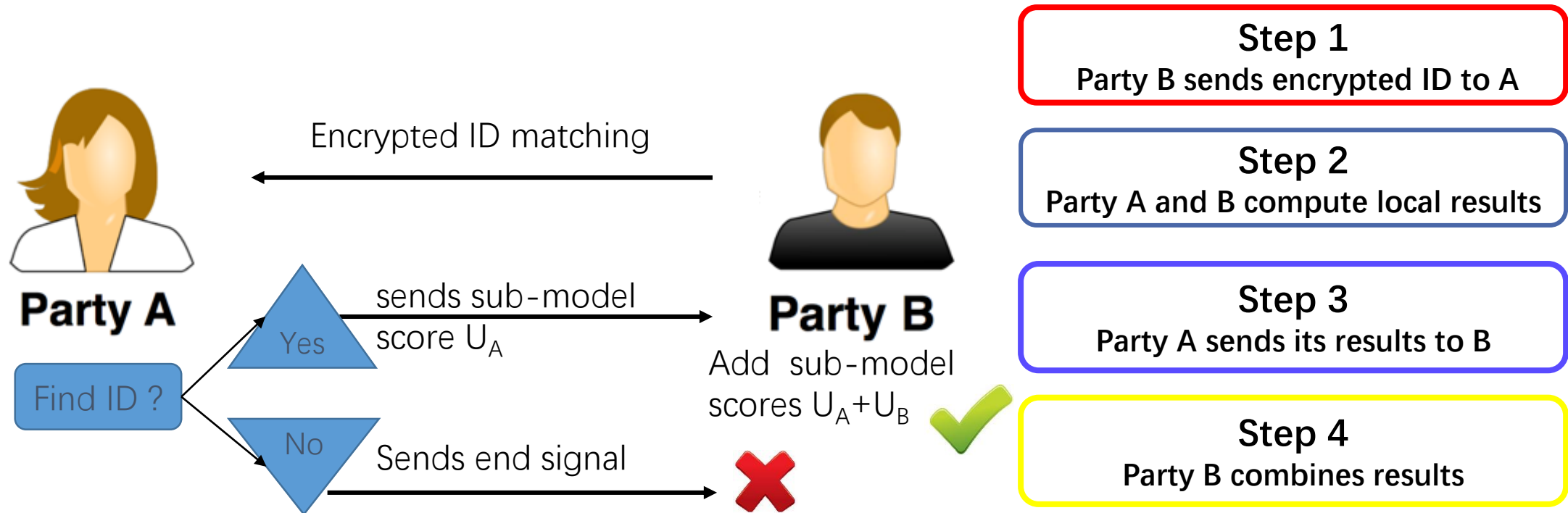
Vertical Federated Learning



Federated Transfer Learning: Concepts and Applications. Qiang Yang, Yang Liu and Tianjian Chen. ACM TIST 2019.

Privacy-Preserving inference

- Suppose a new user ID arrives at Party B,



Security Analysis

- Security against third-party C
 - all C learns are the masked gradients and the randomness and secrecy of the masked matrix are guaranteed
- Security against each other
 - Party A learns its gradient at each step, but this is not enough for A to learn any information from B
 - inability of solving n equations in more than n unknowns
- Security in the semi-honest setting

XGBoost in Federated Learning

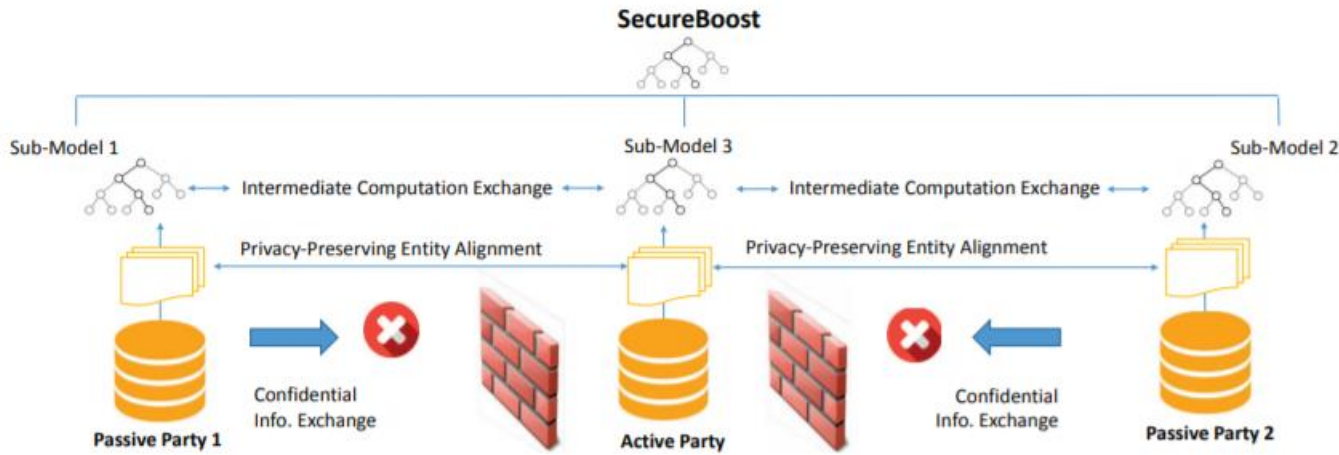
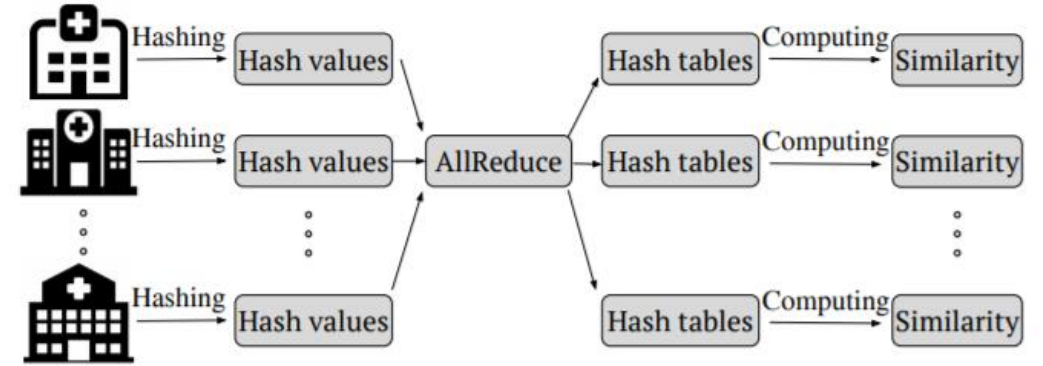


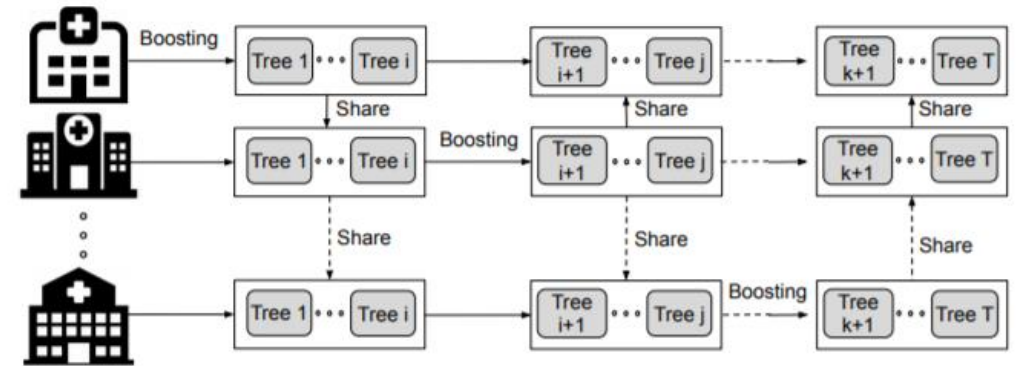
Figure 1: Illustration of the proposed SecureBoost framework

[Kewei Cheng](#), [Tao Fan](#), [Yilun Jin](#), [Yang Liu](#), [Tianjian Chen](#), [Qiang Yang](#), SecureBoost: A Lossless Federated Learning Framework, IEEE Intelligent Systems 2020

GBDT in HFL



(a) The preprocessing stage

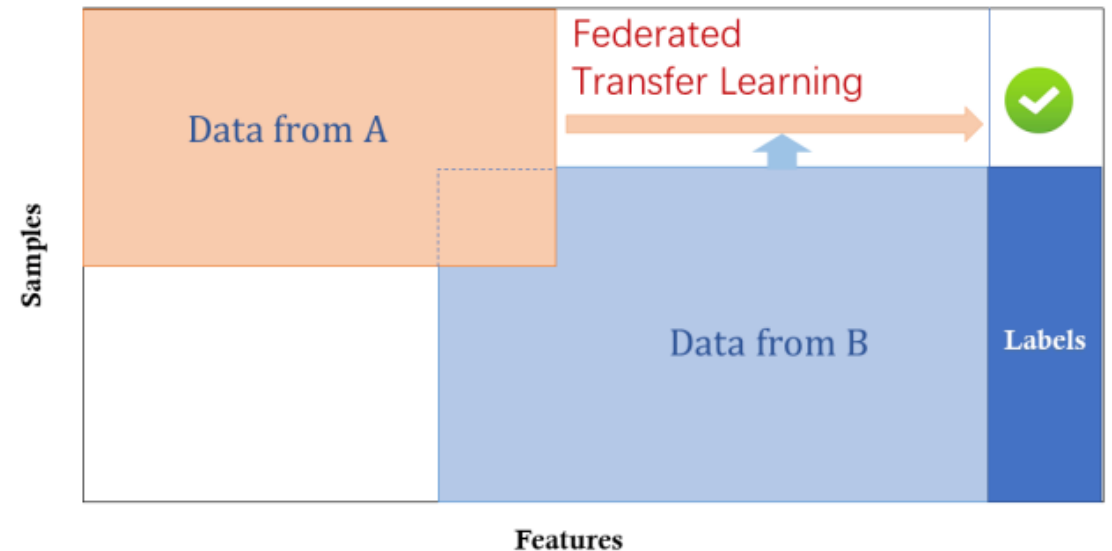
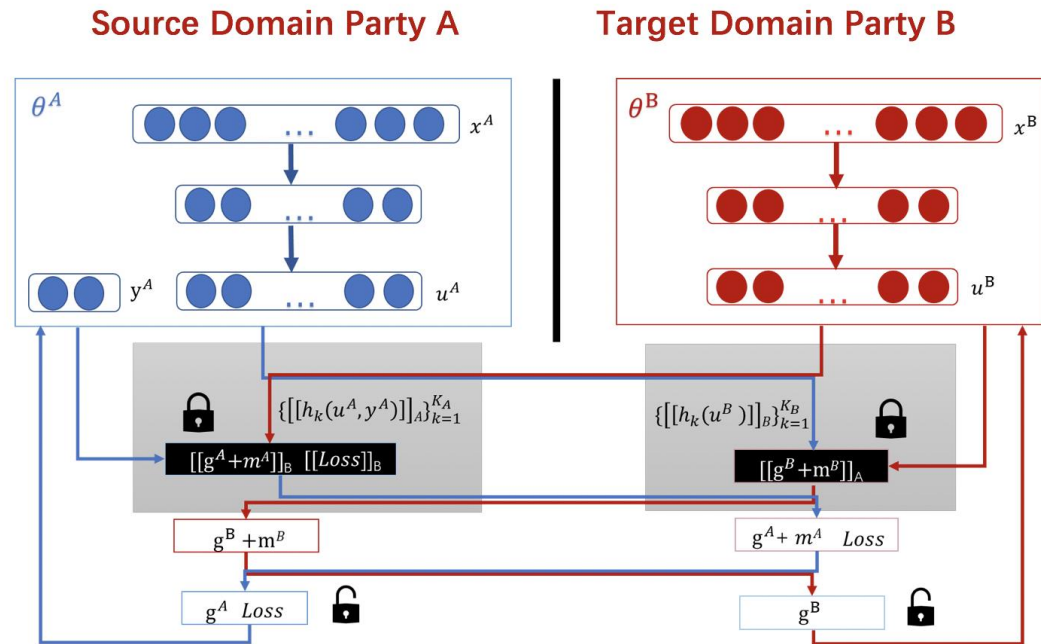


(b) The training stage

[Qinbin Li](#), [Zeyi Wen](#), [Bingsheng He](#), Practical Federated Gradient Boosting Decision Trees, AAAI, 2019

Federated Transfer Learning, FTL

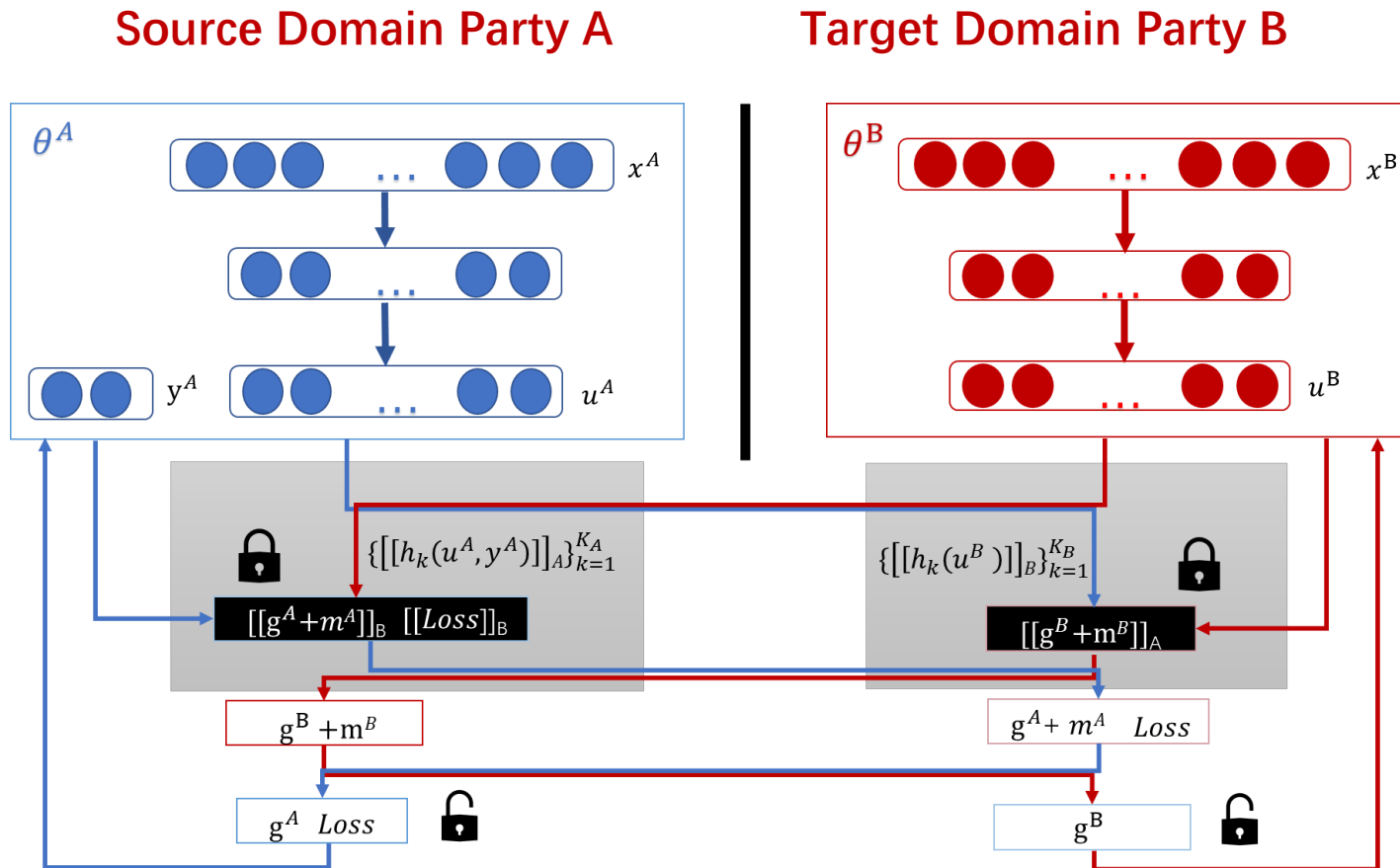
- Parties hold data with different ID and different features (some parties may not have labels).
 - Suitable for federated learning across industries.
 - Common methods include model transfer, instance transfer, feature transfer and domain adaptation, etc.
 - Peer-to-peer architecture is commonly used.
 - **FTL is important to solving the problem of 'small data' and 'unlabeled data'.**



[Yang'19] Qiang Yang, et al., Federated machine learning: Concept and Applications, WeBank, 2019.

[Liu'19] Yang Liu, et al., Secure Federated transfer learning, WeBank, 2018.

Federated Transfer Learning



Step 1
Party A and B send public keys to each other

Step 2
Parties compute, encrypt and exchange intermediate results

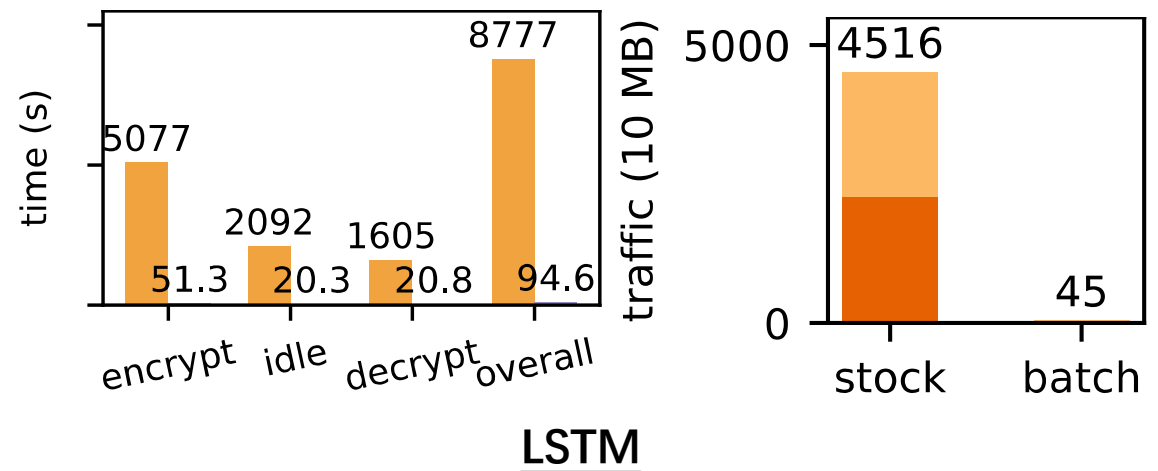
Step 3
Parties compute encrypted gradients, add masks and send to each other

Step 4
Parties decrypt gradients and exchange, unmask and update model locally

Federated Transfer Learning. Yang Liu, Tianjian Chen, Qiang Yang, <https://arxiv.org/pdf/1812.03337.pdf> 2018

Efficiency: BatchCrypt: Efficient Homomorphic Encryption for Cross-Silo Federated Learning

- **Reducing the encryption overhead and data transfer**
 - Quantizing a gradient value into low-bit integer representations
 - Batch encryption: encoding a batch of quantized values to a long integer
- **BatchCrypt is implemented in FATE and is evaluated using popular deep learning models**
 - Accelerating the training by 23x-93x
 - Reducing the netw. footprint by 66x-101x
 - Almost no accuracy loss (<1%)

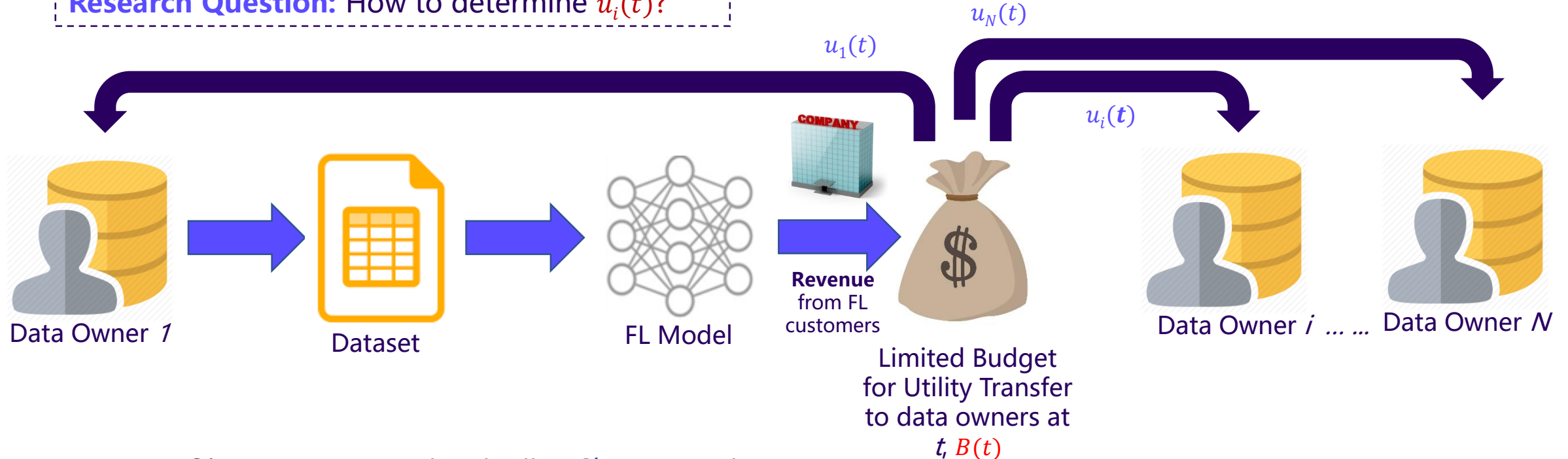


C. Zhang, S. Li, J. Xia, W Wang, F Yan, Y. Liu, BatchCrypt: Efficient Homomorphic Encryption for Cross-Silo Federated Learning, USENIX ATC'20 (accepted)

Incentivize Parties to Join: Federated Learning Exchange

- **Observation:** The success of a federation depends on data owners to share data with the federation
- **Challenge:** How to motivate continued participation by data owners in a federation?

Research Question: How to determine $u_i(t)$?



•Qiang Yang, [Yang Liu](#), [Tianjian Chen](#), [Yongxin Tong](#):

Federated Machine Learning: Concept and Applications. [ACM TIST 10\(2\)](#): 12:1-12:19 (2019)

Advances and Open Problems in Federated Learning

Peter Kairouz^{7*} H. Brendan McMahan^{7*} Brendan Avent²¹ Aurélien Bellet⁹
Mehdi Bennis¹⁹ Arjun Nitin Bhagoji¹³ Keith Bonawitz⁷ Zachary Charles⁷
Graham Cormode²³ Rachel Cummings⁶ Rafael G.L. D'Oliveira¹⁴
Salim El Rouayheb¹⁴ David Evans²² Josh Gardner²⁴ Zachary Garrett⁷
Adrià Gascón⁷ Badih Ghazi⁷ Phillip B. Gibbons² Marco Gruteser^{7,14}
Zaid Harchaoui²⁴ Chaoyang He²¹ Lie He⁴ Zhouyuan Huo²⁰
Ben Hutchinson⁷ Justin Hsu²⁵ Martin Jaggi⁴ Tara Javidi¹⁷ Gauri Joshi²
Mikhail Khodak² Jakub Konečný⁷ Aleksandra Korolova²¹ Farinaz Koushanfar¹⁷
Sanmi Koyejo^{7,18} Tancrede Lepoint⁷ Yang Liu¹² Prateek Mittal¹³
Mehryar Mohri⁷ Richard Nock¹ Ayfer Özgür¹⁵ Rasmus Pagh^{7,10}
Mariana Raykova⁷ Hang Qi⁷ Daniel Ramage⁷ Ramesh Raskar¹¹
Dawn Song¹⁶ Weikang Song⁷ Sebastian U. Stich⁴ Ziteng Sun³
Ananda Theertha Suresh⁷ Florian Tramèr¹⁵ Praneeth Vepakomma¹¹ Jianyu Wang²
Li Xiong⁵ Zheng Xu⁷ Qiang Yang⁸ Felix X. Yu⁷ Han Yu¹² Sen Zhao⁷

¹Australian National University, ²Carnegie Mellon University, ³Cornell University,

⁴École Polytechnique Fédérale de Lausanne, ⁵Emory University, ⁶Georgia Institute of Technology,

⁷Google Research, ⁸Hong Kong University of Science and Technology, ⁹INRIA, ¹⁰IT University of Copenhagen,

¹¹Massachusetts Institute of Technology, ¹²Nanyang Technological University, ¹³Princeton University,

¹⁴Rutgers University, ¹⁵Stanford University, ¹⁶University of California Berkeley,

¹⁷University of California San Diego, ¹⁸University of Illinois Urbana-Champaign, ¹⁹University of Oulu,

²⁰University of Pittsburgh, ²¹University of Southern California, ²²University of Virginia,

²³University of Warwick, ²⁴University of Washington, ²⁵University of Wisconsin–Madison

Applications



Anti money laundering

- ✓ recall improves 15%
- ✓ Audit efficiency improved by over 50%



Internet + banking Risk modeling

- ✓ Performance keeps increasing with respect to the enriched features



Internet + insurance Insurance pricing

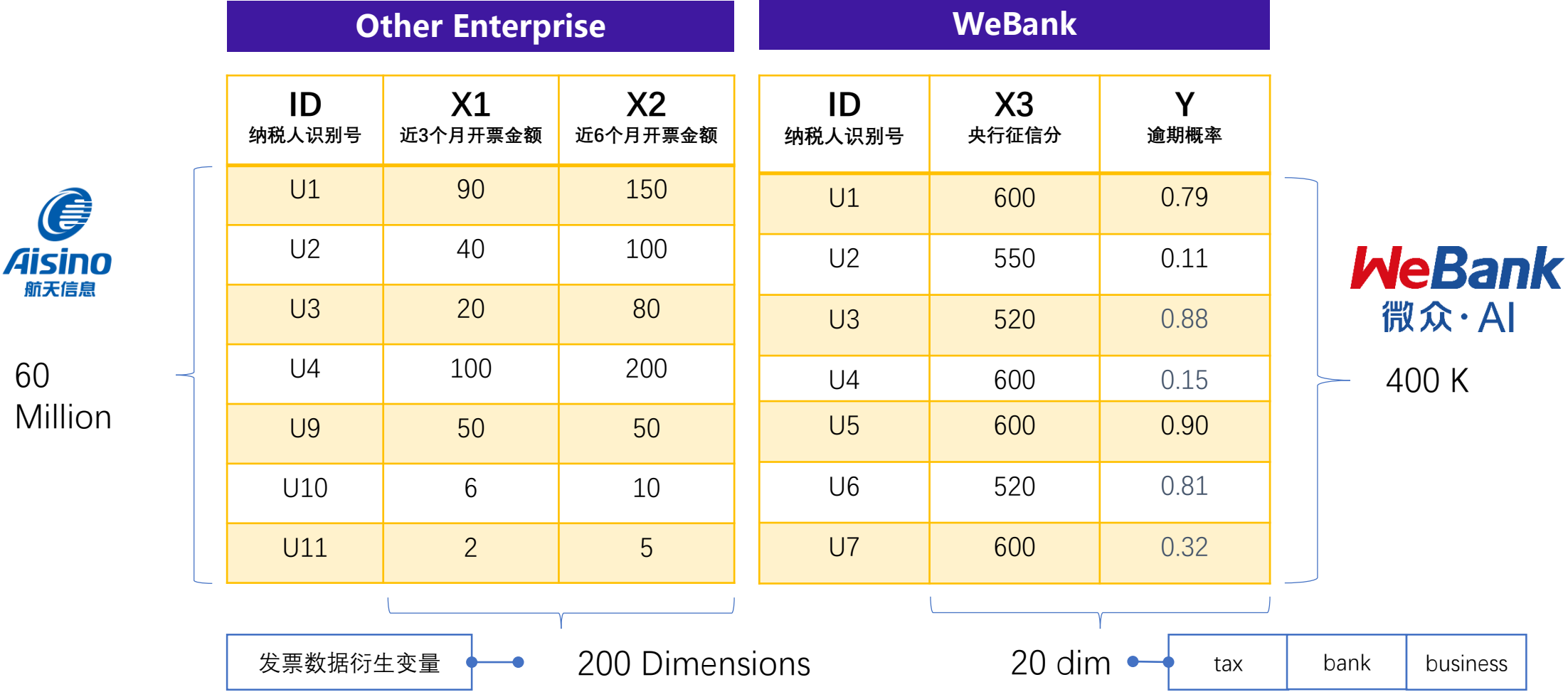
- ✓ Pricing model improves accuracy
- ✓ Coverage ratio is over 90%



Internet + retailers Intelligent marketing

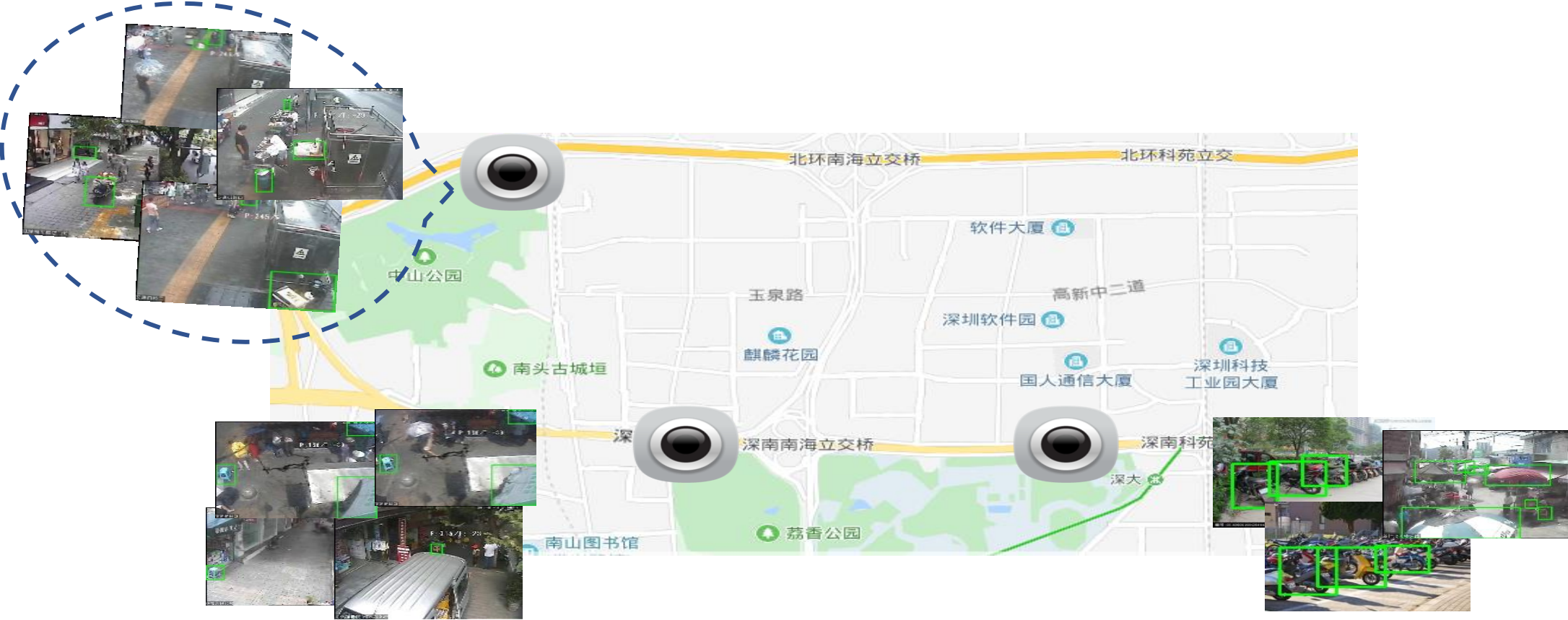
- ✓ Marketing efficiency improves greatly;
- ✓ Better user profile and targeting;

Risk Management with Federated Learning

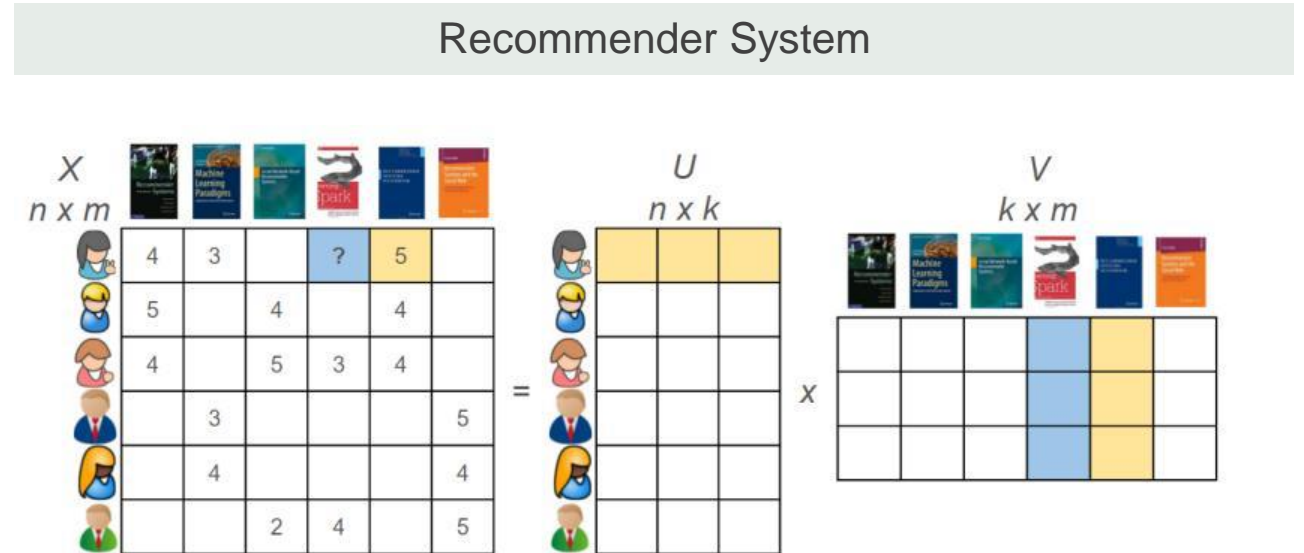
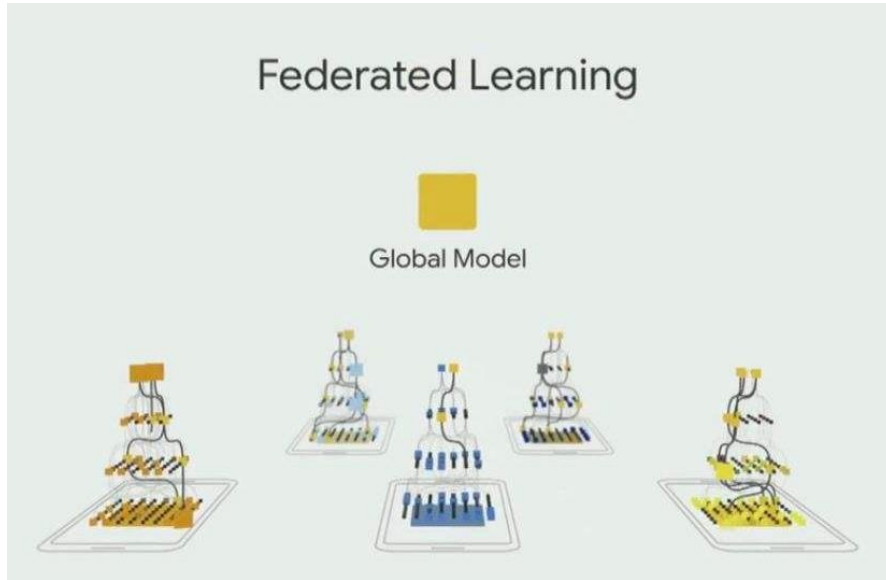


Construction-Site Safety w/ Federated Computer Vision

Webank AI X Extreme Vision



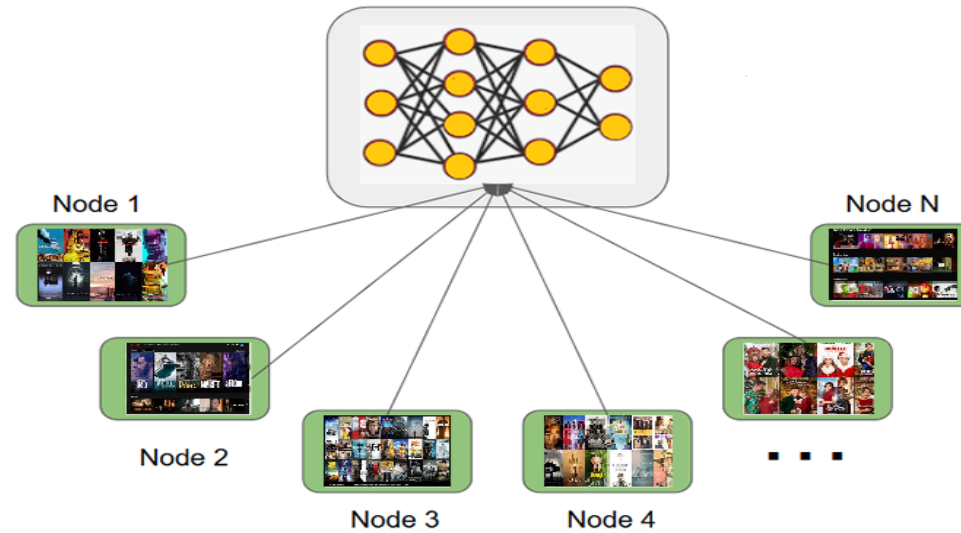
Federated Recommendation



Assumption: a trustworthy 3rd-party as coordinator, which can be removed

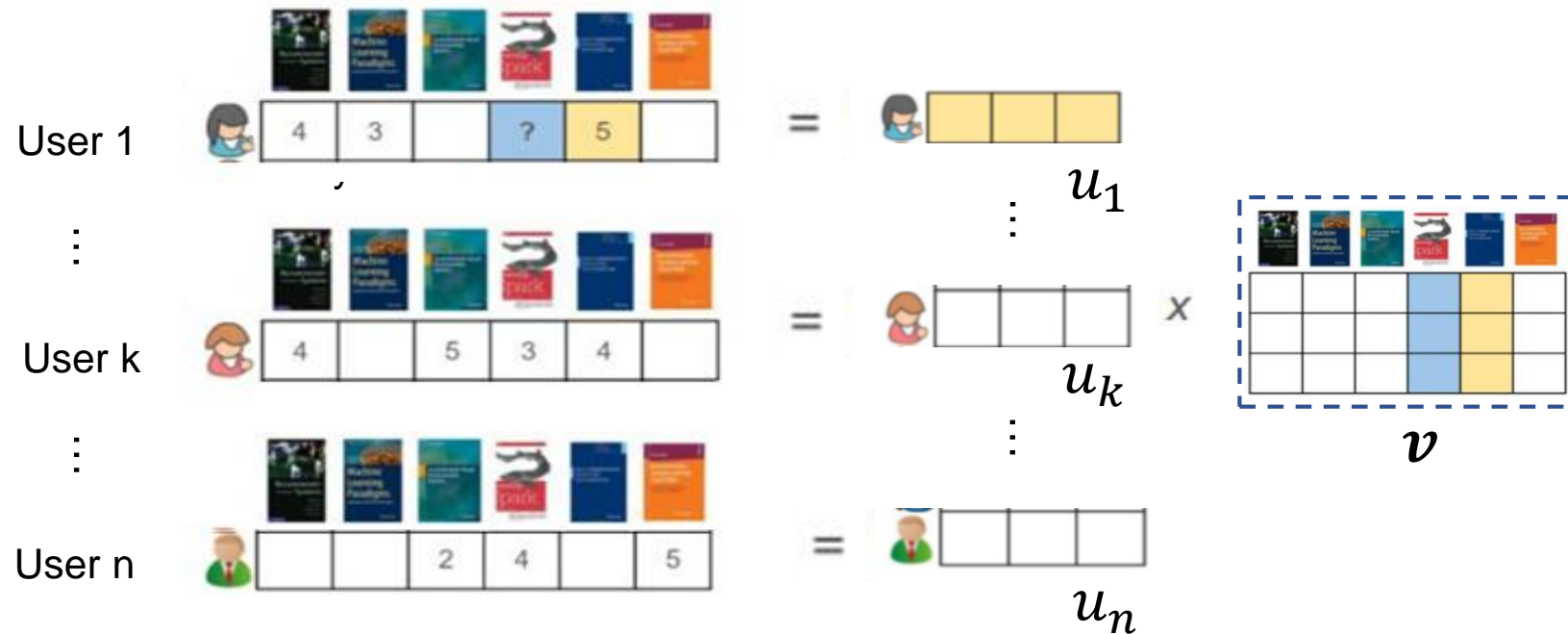
Horizontal Federated Recommendation

Example: movie recommendation with data from individual users



Federated Collaborative Filtering [Ammad et al. 2019]

Intuition: decentralized matrix factorization, each user profile is updated locally,
item profiles are aggregated and updated by coordinator



Loss function
$$\min_{U, V} \frac{1}{M} (r_{i,j} - \langle u_i, v_j \rangle)^2 + \lambda \|U\|_2^2 + \mu \|V\|_2^2$$

Update function

$$\begin{array}{l}
 \boxed{u_i^t = u_i^{t-1} - \gamma \nabla_{u_i} F(U^{t-1}, V^{t-1})} \rightarrow \text{User local updates} \\
 \boxed{v_i^t = v_i^{t-1} - \gamma \nabla_{v_i} F(U^{t-1}, V^{t-1})} \leftarrow \text{Server updates} \rightarrow \text{Gradients from users}
 \end{array}$$

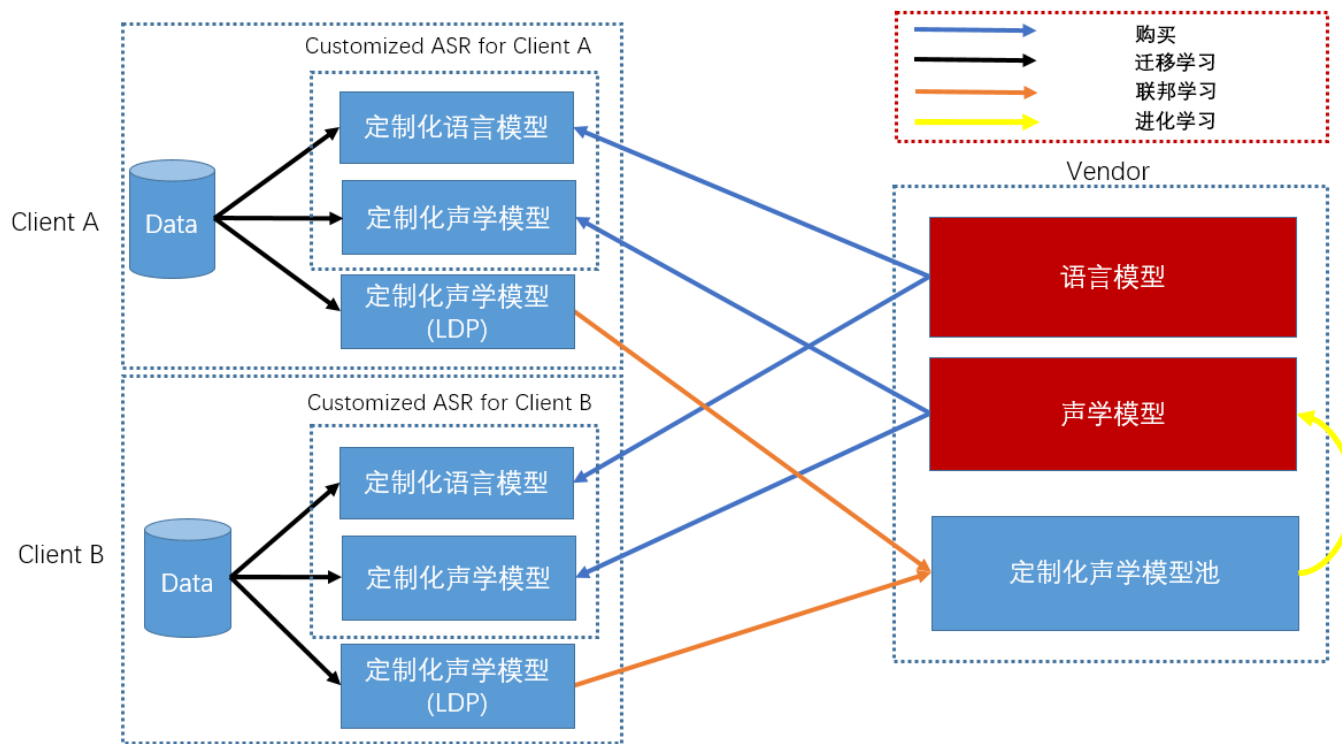
Voice Recognition using Multiple Data Sources



传统语音识别

基于TFE的语音识别

TFE= Transfer + Federated + Evolutionary



Reduced error by 10%~20%

Federated Learning Eco-Systems in B2B Markets: at WeBank

First Industrial Strength Enterprise Federated Learning System

稳定运行2年以上, 数十家客户落地运行

100 Million Level Data Processing, 2000+ Features

Credit Rating in Banking : 12 scenarios

Awards

AAAI 2020 Industry Award

SaaS Queries

40 Million

支持亿级数据隐私保护交集算法,
支撑 20+ 联邦学习算法生产应用,
超**4000万**在线调用

Eco systems

530+ organizations

同比上涨96%, 覆盖 370 家
企业机构, 164 所高校

Patents

74 Approved

In applications: > 400,

Standards

5项标准制定

其中2项国际标准, 3项
国内标准

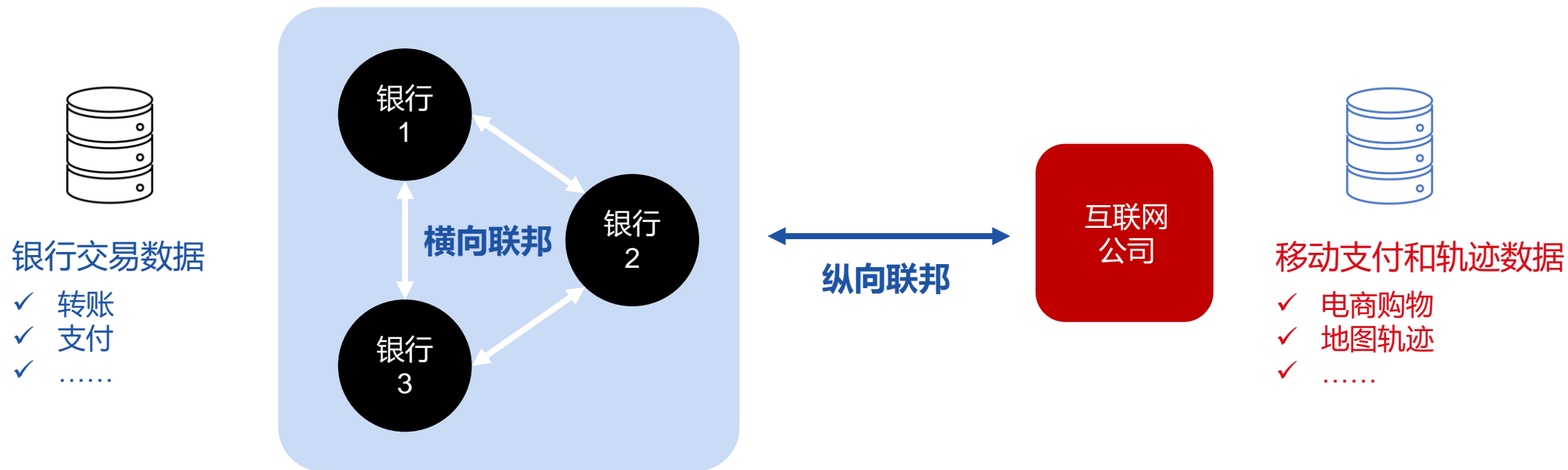
Textbooks

《**联邦学习**》



联合中国银联、鹏城实验室、
平安科技等共同发布了《联邦
学习白皮书2.0》

Anti-money Laundering



通过横向联邦扩充反洗钱样本，构建基础反洗钱模型 → 通过纵向联邦扩充客户特征维度，进一步优化模型效果

Multiple Data Sources in Credit Rating

多方数据本地建模

联邦建模

联邦评分



安全联合多方多维数据，提高AI模型精度

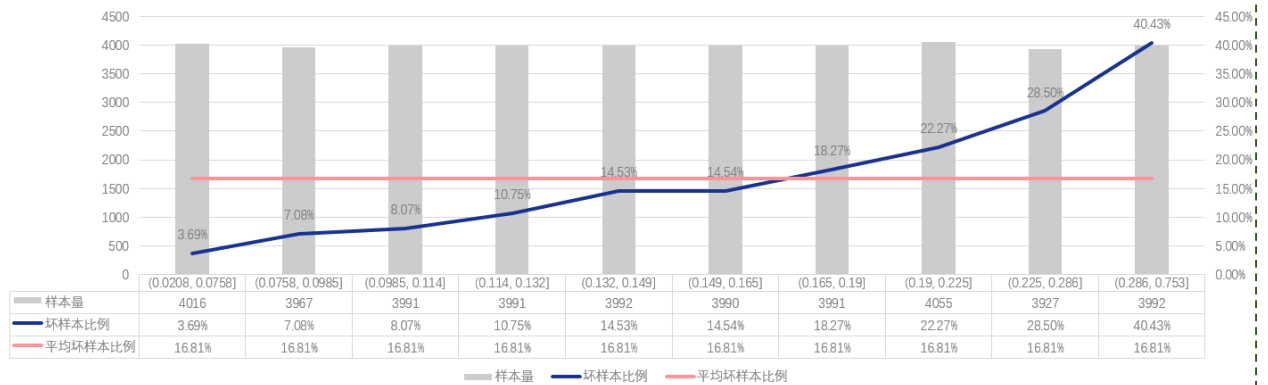
分数对应资质好坏

反欺诈/贷前评分/贷后监测

反欺诈评分举例

模型效果: $AUC=0.70$ $KS=30$, 尾部分组坏样本比例是平均样本比例的2.4倍

使用方式: 1) 单一策略; 2) 入模变量; 3) 决策矩阵



Insurance Industry

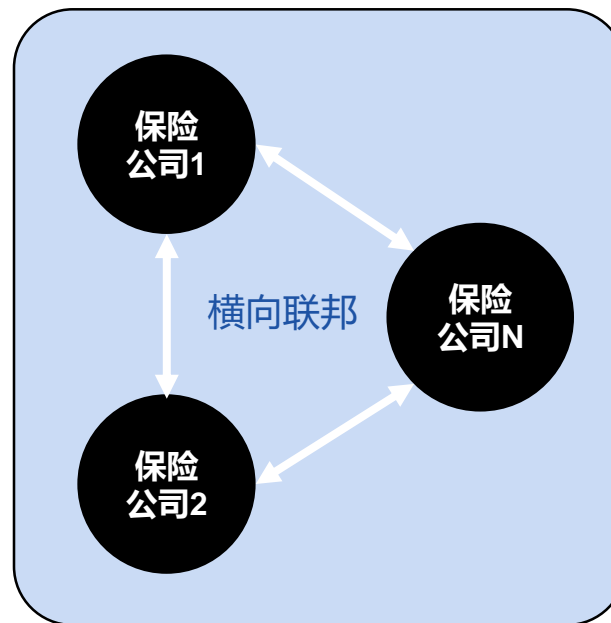
WeBank and Swiss Re



协助再保公司建立承保人（保险公司）的车险索赔概率模型：纵向联邦引入和挖掘互联网大数据“从人因子”，横向联邦扩大承保人传统因子数据集规模，从而实现对车主进行精准画像和风险分析

Internet Data

- ✓ 出行数据
- ✓ 消费数据
- ✓ 信息偏好
- ✓ 车辆违章数据
- ✓



Insurance

- ✓ 承保数据
- ✓ 理赔数据
- ✓ 车联网数据
- ✓

Federated Recommendation

Example: movie and book recommendation with data from **two different data sources**



豆瓣读书

	4	3		?	5	
	5		4		4	
	4		5	3	4	
		3				5
		4				4
			2	4		5



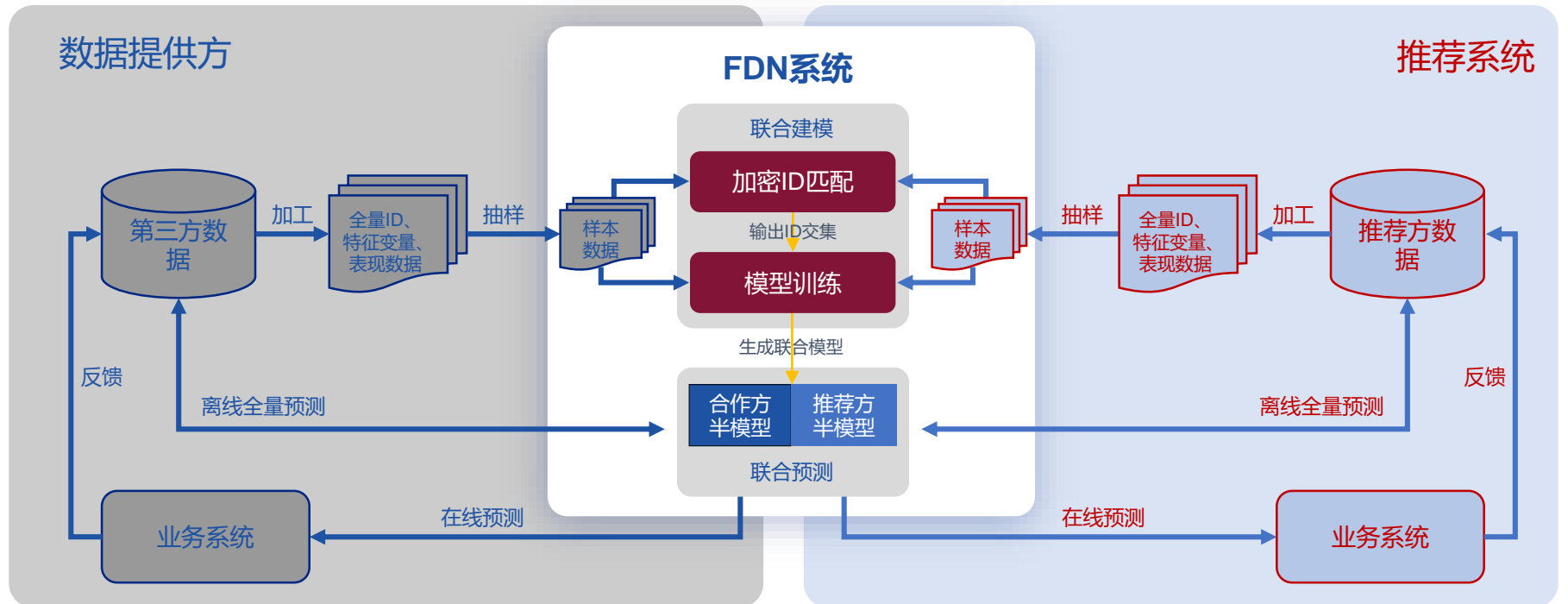
No data exchange

		3			4	4
					5	3
		4				
				4		2
	5					
	3	2	4			

Party A

Party B

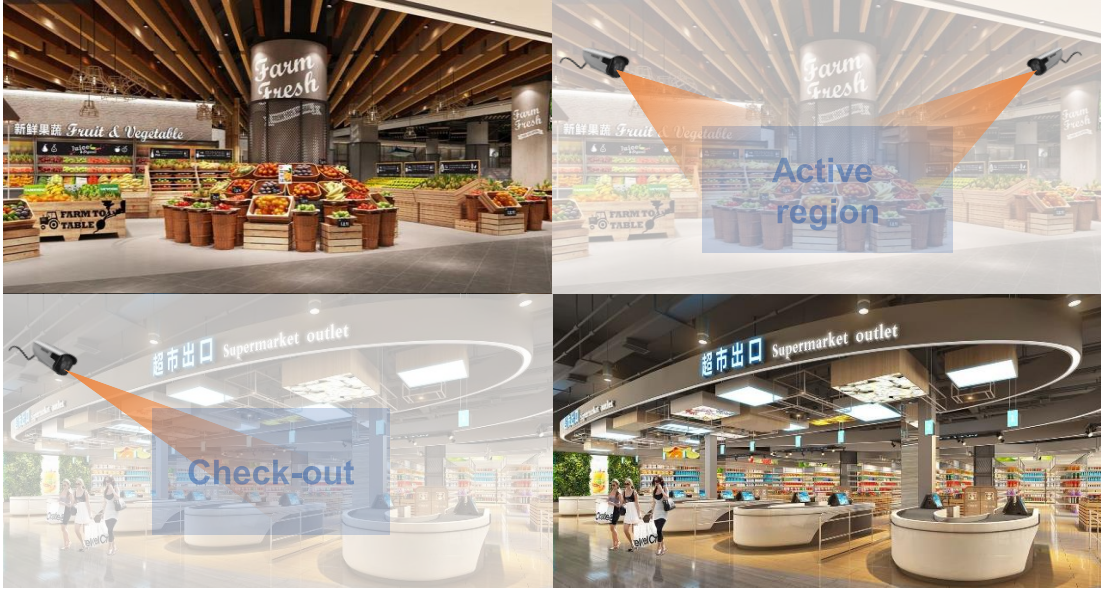
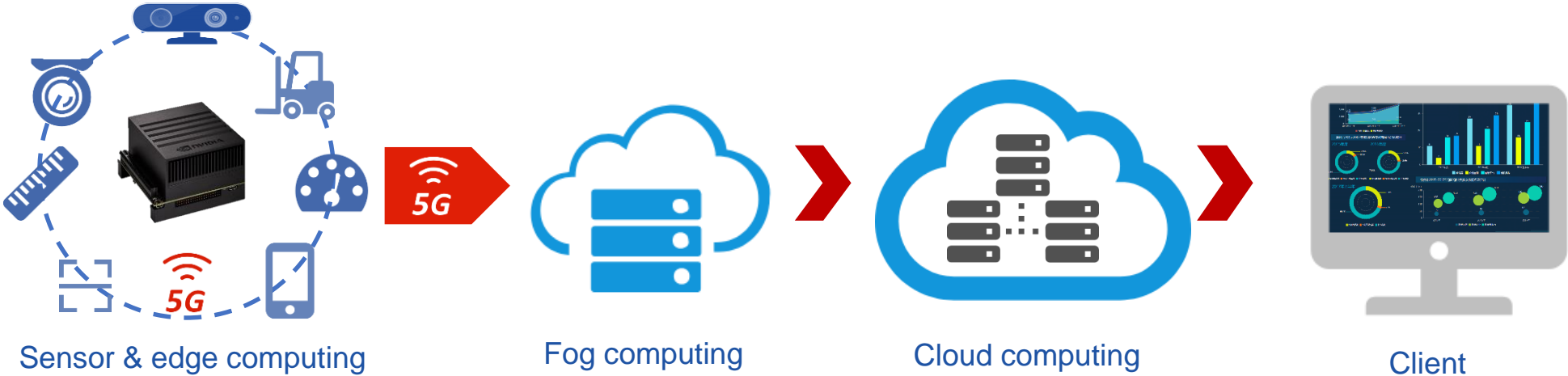
联合建模、预测示意图
—— 安全合规的数据合作
过程



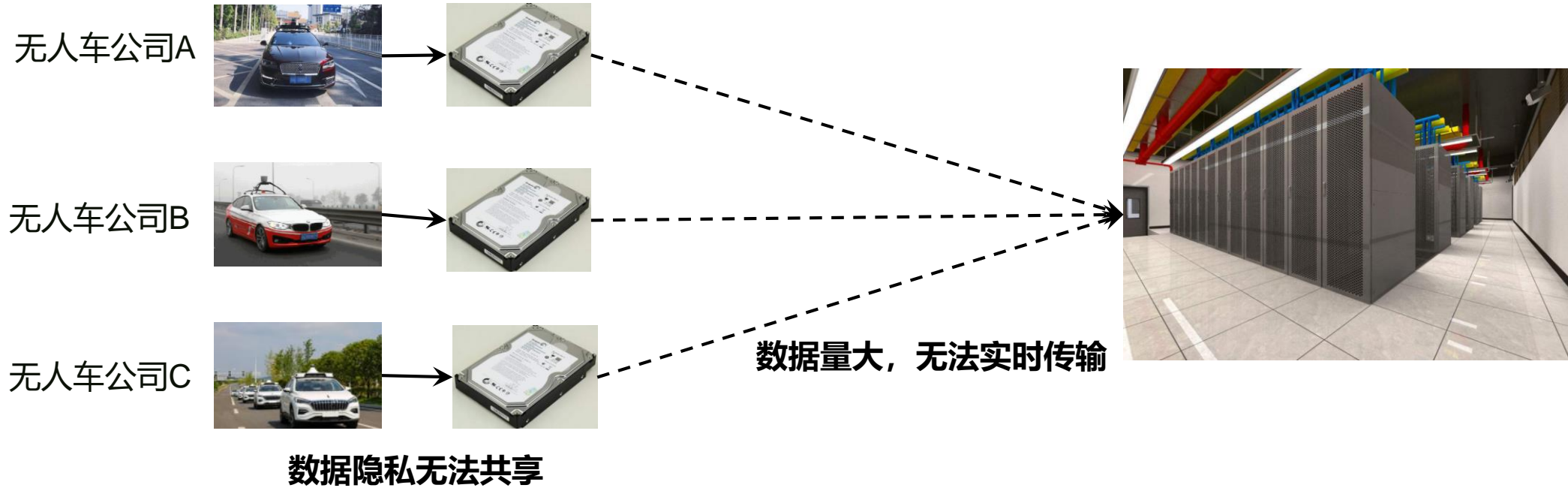
注: 客户ID包括但不限于客户身份证号码、手机号、设备ID (imei) 等; 联合建模过程由拥有Y (表现数据) 的一方发起;

IoT Applications

Perception engine
analysis engine
recognition engine



FedEdge: Federated Edge Computing



Computer Vision using Federated Learning

装备制造业、物联网AIOT、智慧安防等行业，依托联邦学习，进行视觉市场的场景拓宽

优势：

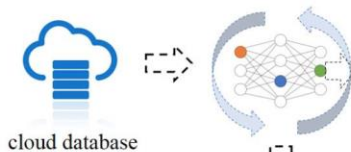
- 相对于本地建模进一步提升算法准确率
- 形成网络效应，降低长尾应用成本，提升视觉业务总体利润率

FedVision –由联邦学习提供支持的在线视觉对象检测平台

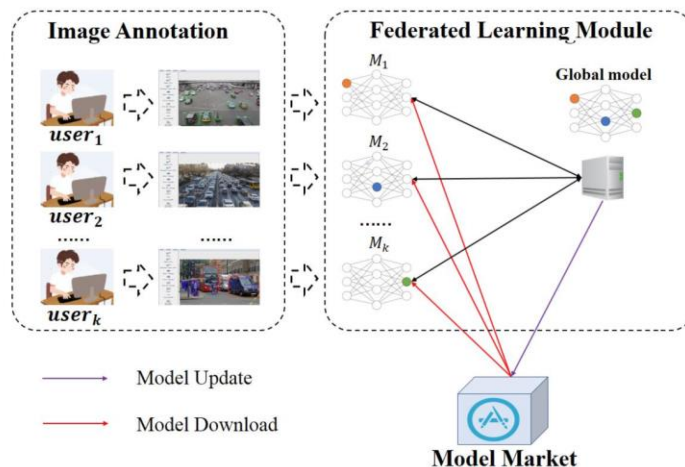
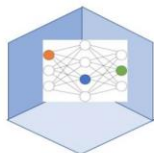
1. Image Annotation



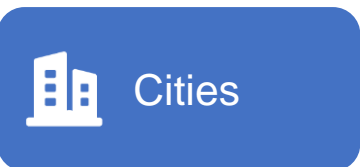
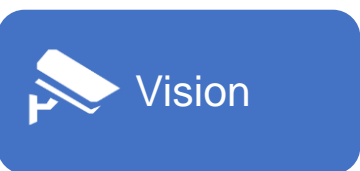
2. Centralized training



3. Online inference



Tasks:
行人检测
出行检测
区域检测
设备异常检测
安全帽检测
火焰检测
烟雾检测
.....



Local models are limited

Federated Vision improves accuracy by 15%

IEEE International Standard for Federated Learning

<https://sagroups.ieee.org/3652-1/>



More Standards

China Communications Standards Association

《基于联邦学习的数据流通产品 技术要求与测试方法
Data circulation products based on federated learning: Technical requirements and testing methods》

IEEE International Standard for Federated Learning

IEEE P3652.1

《Guide for Architectural Framework and Application of Federated Machine Learning(联邦学习基础架构与应用)》

- **March 2021**, The first international federated learning standard
- 20+ members
- 6 working group meetings
- 10 kinds of federated learning application scenarios specifications



IEEE SA
STANDARDS STORE

SHOP by Category Search IEEE Standards MY

IEEE 3652.1-2020
IEEE Guide for Architectural Framework and Application of Federated Machine Learning
STANDARD by IEEE, 03/19/2021
View all product details

WeBank 微众银行 IEEE Baidu 百度
创新工场 INNOVATION VENTURES Hisense Paradigm 第四范式
CLUSTAR 星云 腾讯云 京东商城
中国移动 China Mobile eduworks

WeBank AI: FATE - Federated Learning Open Source Platform

 The AAAI Conference
on Artificial Intelligence

2月
微众AI 领衔推动联邦学
习国际标准制定



1月
微众AI亮相 AAAI 会议
发布联邦学习开源框架
FATE0.1版本



5月
主要合作伙伴切换到FATE
微众与瑞士再保险签订战略合
作协议，推动联邦学习在再保
险业应用



7月 微众AI牵头的国内
首个联邦学习标准正式
出台



8月 微众AI主导首个国
际联邦学习学术研讨会


T/AIOSS
中国人工智能开源软件发展联盟标准
AIOSS-03-2019

6月 微众联邦学习开源
项目加入Linux基金会

Multi-Center Medical Modeling using Federated Learning

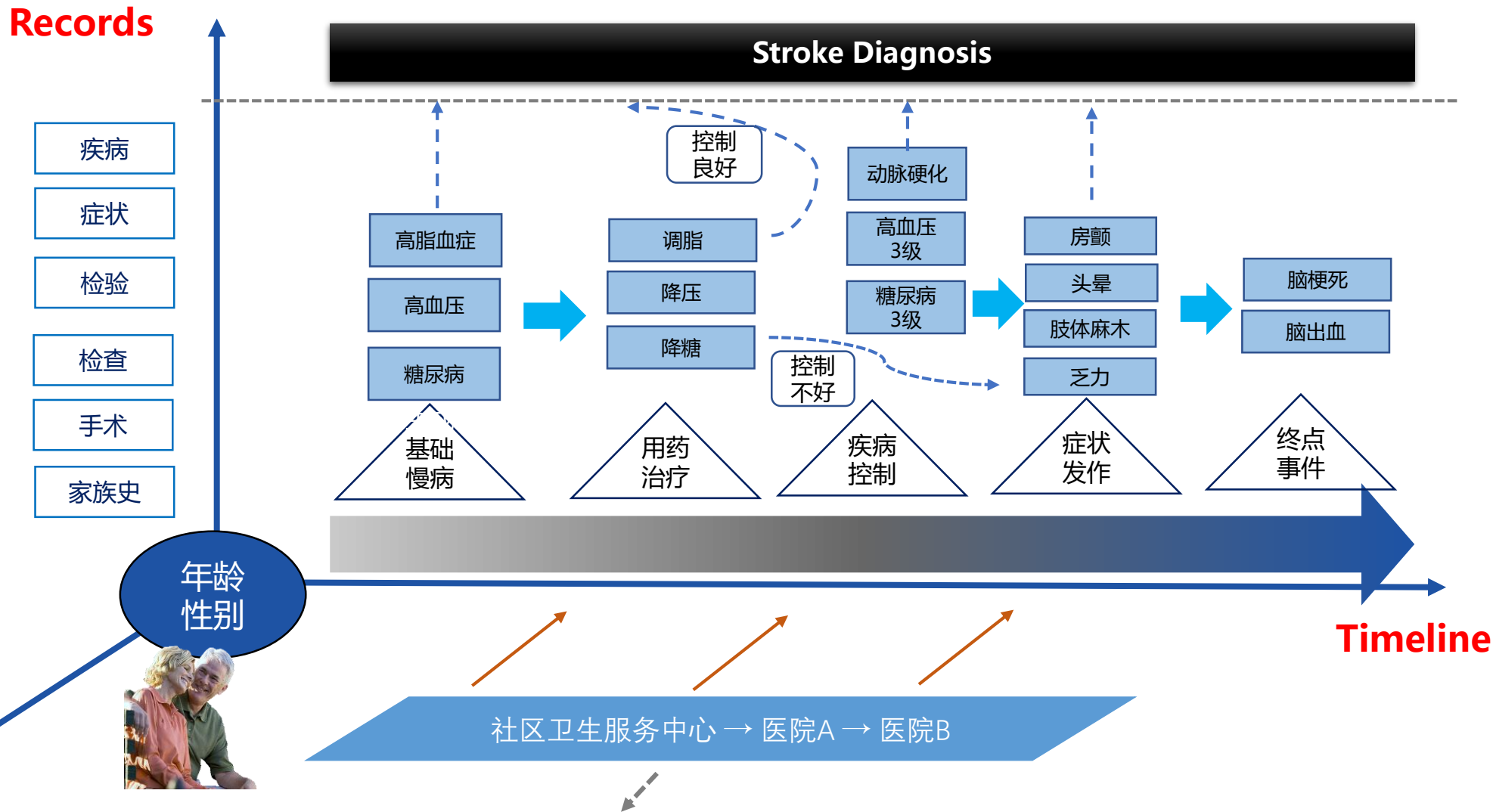
病例：58岁女性，有冠心病、高血压和多次头痛乏力病史、2017年11月30号因头晕四肢麻木被误诊颈椎病和眩晕综合征，模型预测卒中风险78%，最终于2018年5月15号突发脑梗死住院



- Predication based on Time Series, Unstructured Data

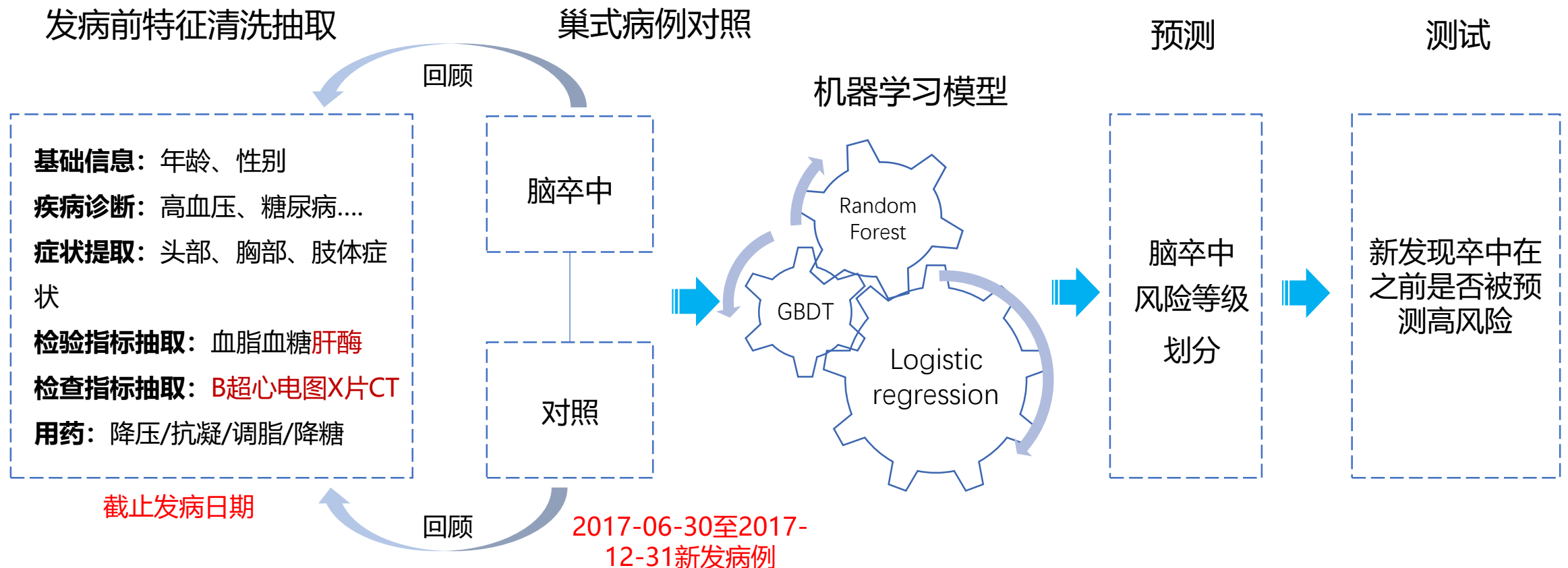
Stroke Patients: How to Predict based on Symptoms?

Using Data from Multiple Diagnostic Centers



Stroke Patients: How to Predict based on Symptoms?

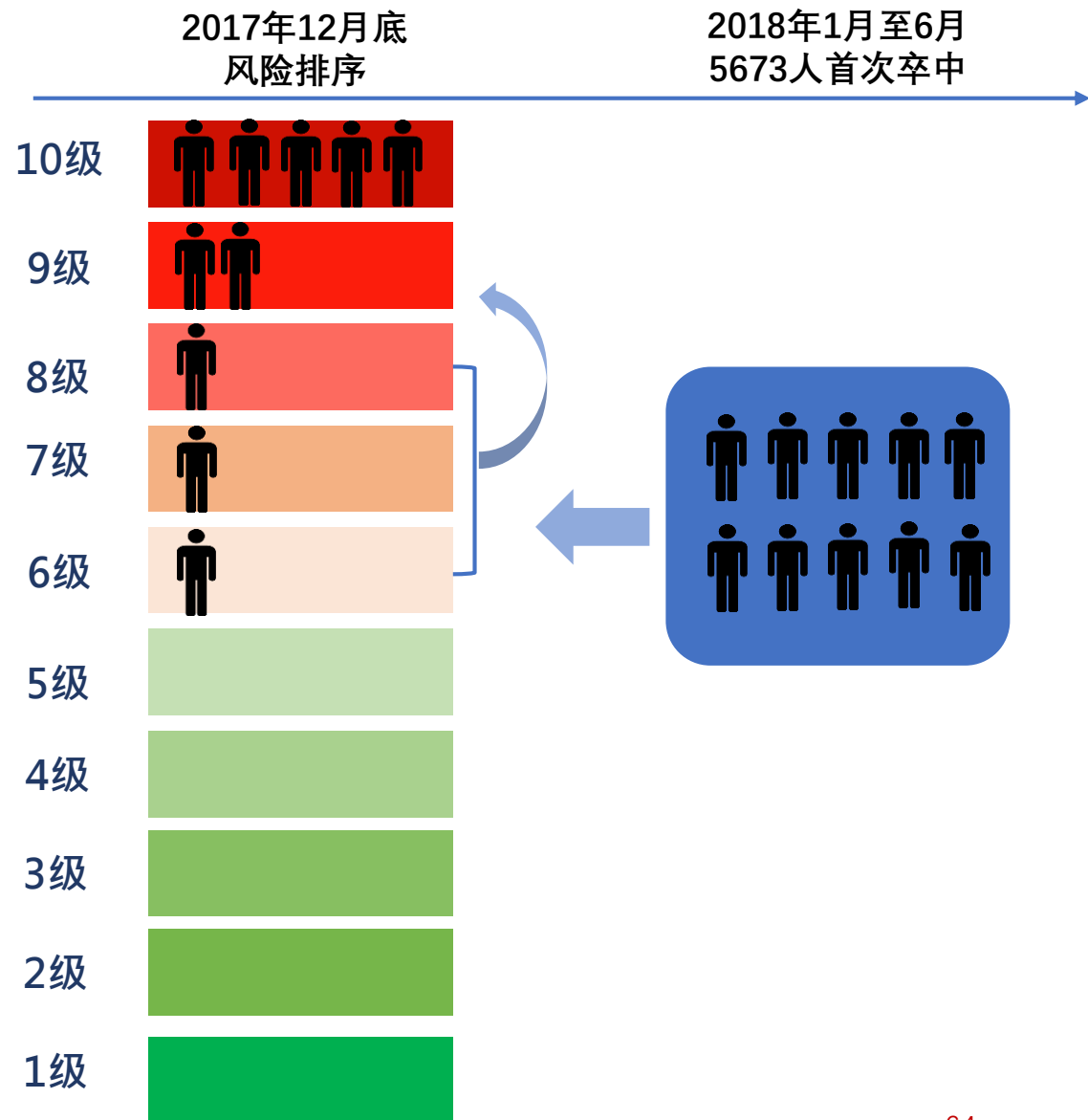
- 抽取2017-06-30至2017-12-31新发现脑卒中（包括脑梗、短暂性脑缺血发作、脑动供血不足、非创伤和血管畸形的脑出血和蛛网膜下腔出血）病例，筛选在发病前至少有两次就诊记录者6502人
- 建立巢式病例对照研究队列：随机抽取与卒中就诊时间相同的其它就诊病例作为对照



Stroke Patients: How to Predict based on Symptoms?

Prediction Results:

- 根据脑卒中发病风险概率模型在2017年12月底对全体人群进行预测和风险排序，设置10个风险等级
- 其中5673人在2018年1-6月期间首次发生脑卒中，每10个发病者中有5人处于10级风险，2人处于9级风险，1人处于8级风险，2人处于6-7级风险，0人处于1-5级风险



Stroke Patients: How to Predict based on Symptoms?

Hubei Hospitals

现状与挑战：医疗大数据模型面临数据孤岛难题

解决方案：联邦学习

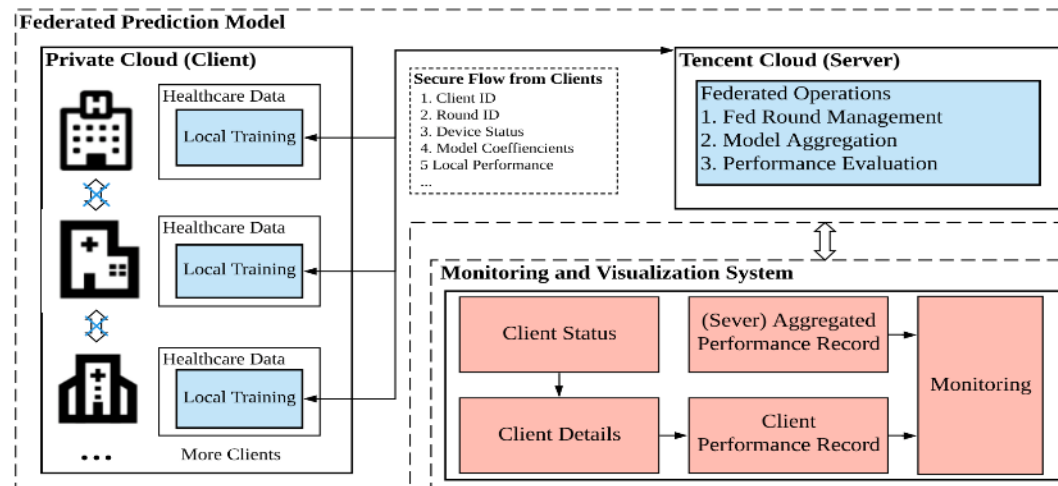
成果：

- 完成多节点部署模型效果POC验证
- 提升了两家小医院D和E的模型AUC 10%和20%以上
- 共建联合实验室

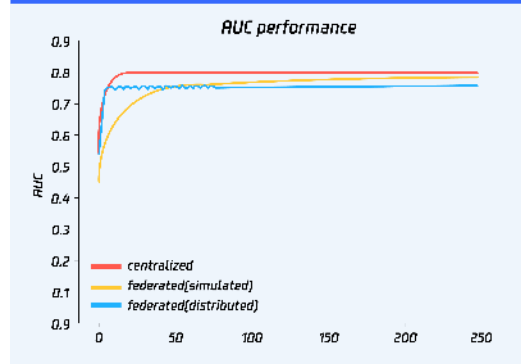
某市TOP 5医院真实数据分布情况

医院ID	样品总数	正样品数	负样品数	特征维度
A Hospital	132,631	4,271	128,360	119
B Hospital	36,118	1,097	35,021	119
C Hospital	18,876	1,196	17,680	119
D Hospital	17,123	100	17,023	119
E Hospital	11,076	68	11,008	119

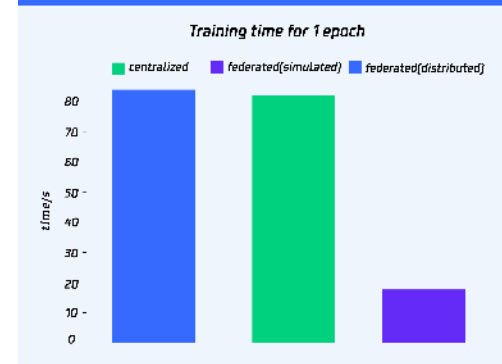
Federated Learning



联邦学习各模型AUC表现对比



联邦学习训练时间对比



Ce Ju et al., "The Privacy-preserving Technology to Help Millions of People: Federated Prediction Model for the Risk of Stroke," IJCAI Workshop on Federated Learning, 2020

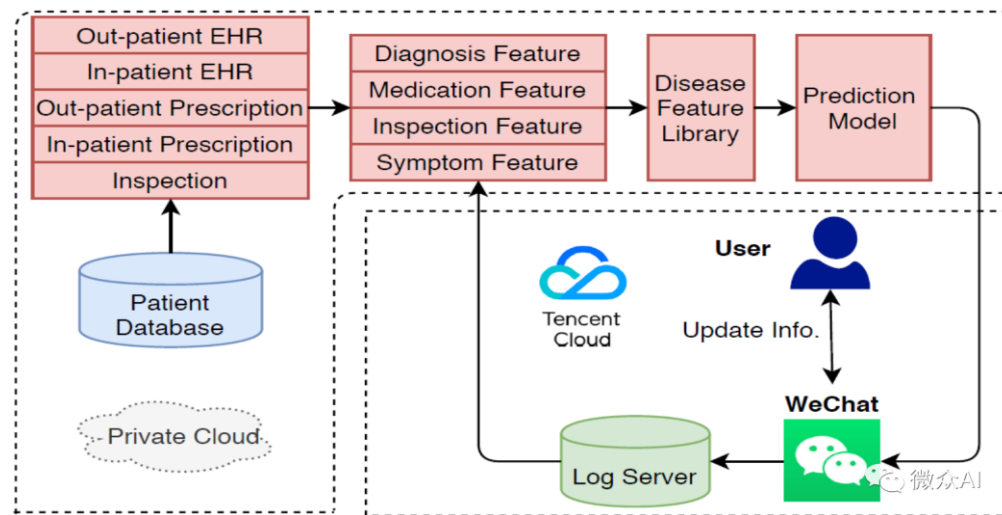
Webank+Tencent Joint Team

- Accuracy: Improvement over 80%,
- Small hospitals: 10-20%

研究论文被FL-IJCAI'20收录

Privacy-Preserving Technology to Help Millions of People: Federated Prediction Model for Stroke Prevention

Ce Ju^{1,*}, Ruihui Zhao^{2,*}, Jichao Sun^{2,*}, Xiguang Wei^{1,*},
Bo Zhao², Yang Liu¹, Hongshan Li³, Tianjian Chen¹,
Xinwei Zhang⁴, Dashan Gao^{5,6}, Ben Tan¹, Han Yu⁷ and Yuan Jin⁸



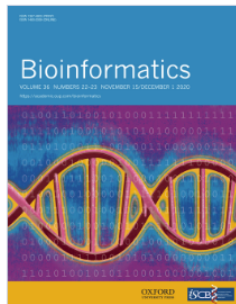
FL-QSAR

An federated learning-based QSAR prototype for collaborative drug discovery

Bioinformatics



Issues Advance articles Submit ▾ Purchase Alerts About ▾



Volume 36, Issue 22-23
1 December 2020

FL-QSAR: a federated learning-based QSAR prototype for collaborative drug discovery

Shaoqi Chen, Dongyu Xue, Guohui Chuai, Qiang Yang ✉, Qi Liu ✉

Bioinformatics, Volume 36, Issue 22-23, 1 December 2020, Pages 5492–5498,
<https://doi.org/10.1093/bioinformatics/btaa1006>

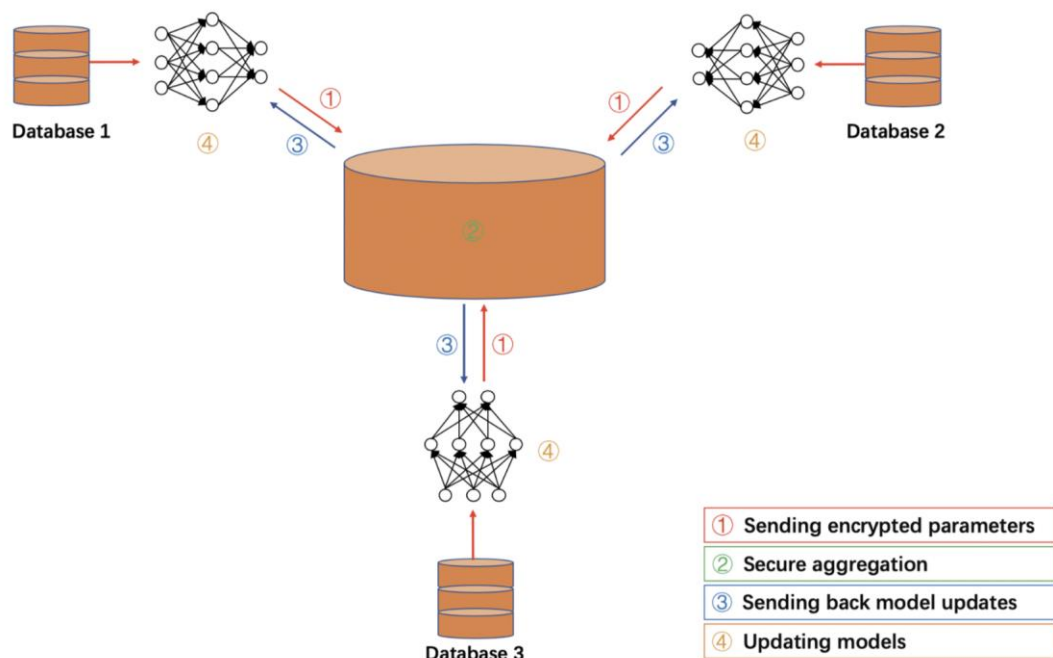
Published: 08 December 2020 **Article history** ▾

FL-QSAR

Federated Learning based New Drug Structure Discovery:

Accuracy same as putting the data all together.

FL Based platform: FL-QSAR



4-step training process:

- (1) 各制药机构方发送加密模型参数到聚合中心
- (2) 聚合中心进行参数聚合
- (3) 聚合中心将聚合后的参数发回各制药机构用户
- (4) 各制药机构用户使用聚合后的参数更新模型

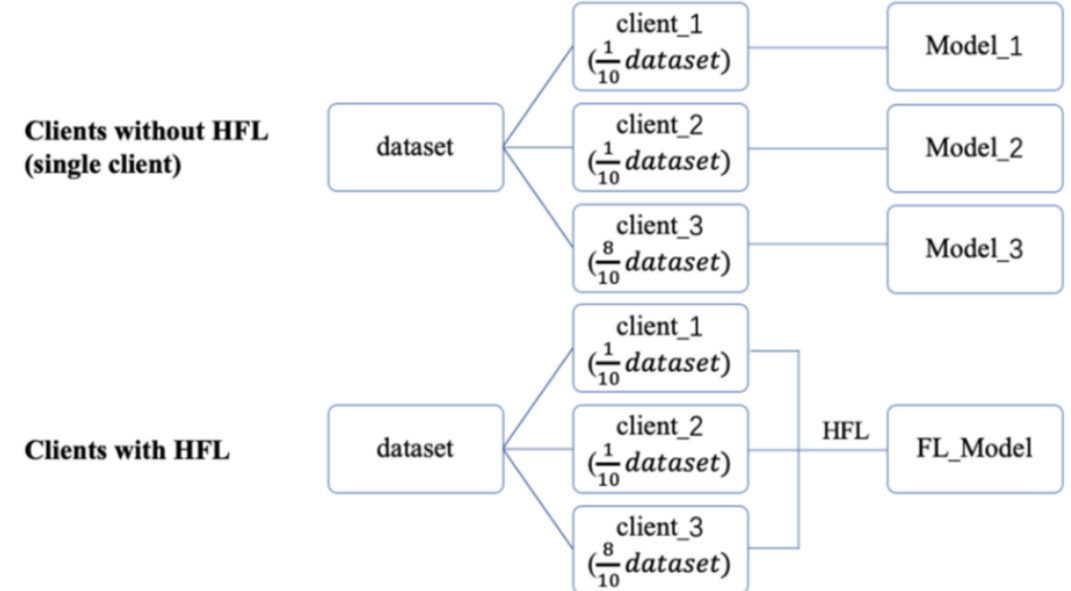
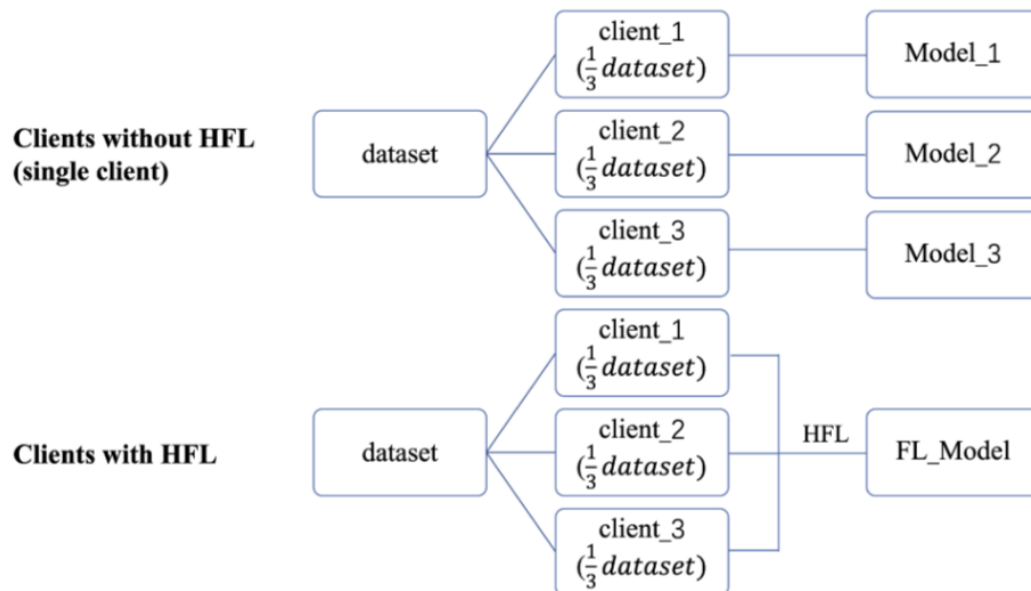
Shaoqi Chen et al. Qiang Yang, Qi Liu, 2020, *Bioinformatics*.

FL-QSAR Experiments

14 QSAR data sets

Two or more distributed collaboration parties

Data imbalanced: three centers



Shaoqi Chen et al. Qiang Yang, Qi Liu, 2020, *Bioinformatics*.

Benchmark QSAR Data

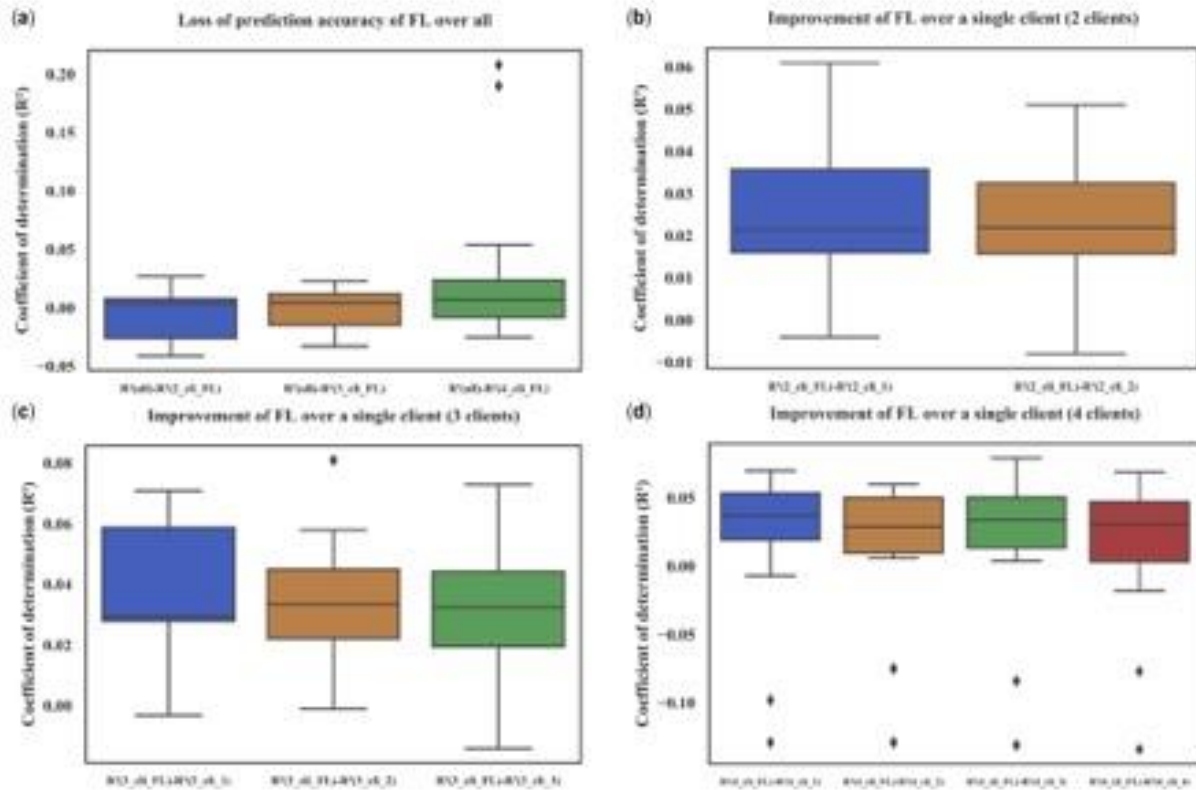
Table 1. Benchmark datasets tested in FL-QSAR

Datasets	Data type	Description	Number of molecules	Number of feature descriptors
3A4	ADME	CYP P450 3A4 inhibition $-\log(\text{IC}_{50})$ M	50 000	9491
CB1	Target	Binding to cannabinoid receptor 1 $-\log(\text{IC}_{50})$ M	11 640	5877
DPP4	Target	Inhibition of dipeptidyl peptidase 4 $-\log(\text{IC}_{50})$ M	8327	5203
HIVINT	Target	Inhibition of HIV integrase in a cell based assay $-\log(\text{IC}_{50})$ M	2421	4306
HIVPROT	Target	Inhibition of HIV protease $-\log(\text{IC}_{50})$ M	4311	6274
LOGD	ADME	LogD measured by HPLC method	50 000	8921
METAB	ADME	Percent remaining after 30 min microsomal incubation	2092	4595
NK1	Target	Inhibition of neurokinin1 (substance P) receptor binding $-\log(\text{IC}_{50})$ M	13 482	5803
OX1	Target	Inhibition of orexin 1 receptor $-\log(\text{K}_i)$ M	7135	4730
OX2	Target	Inhibition of orexin 2 receptor $-\log(\text{K}_i)$ M	14 875	5790
PGP	ADME	Transport by p-glycoprotein $\log(\text{BA}/\text{AB})$	8603	5135
PPB	ADME	Human plasma protein binding $\log(\text{bound}/\text{unbound})$	11 622	5470
RAT_F	ADME	$\log(\text{rat bioavailability})$ at 2 mg/kg	7821	5698
TDI	ADME	Time dependent 3A4 inhibitions $\log(\text{IC}_{50}$ without NADPH/ IC_{50} with NADPH)	5559	5945
THROMBIN	Target	Human thrombin inhibition $-\log(\text{IC}_{50})$ M	6924	5552

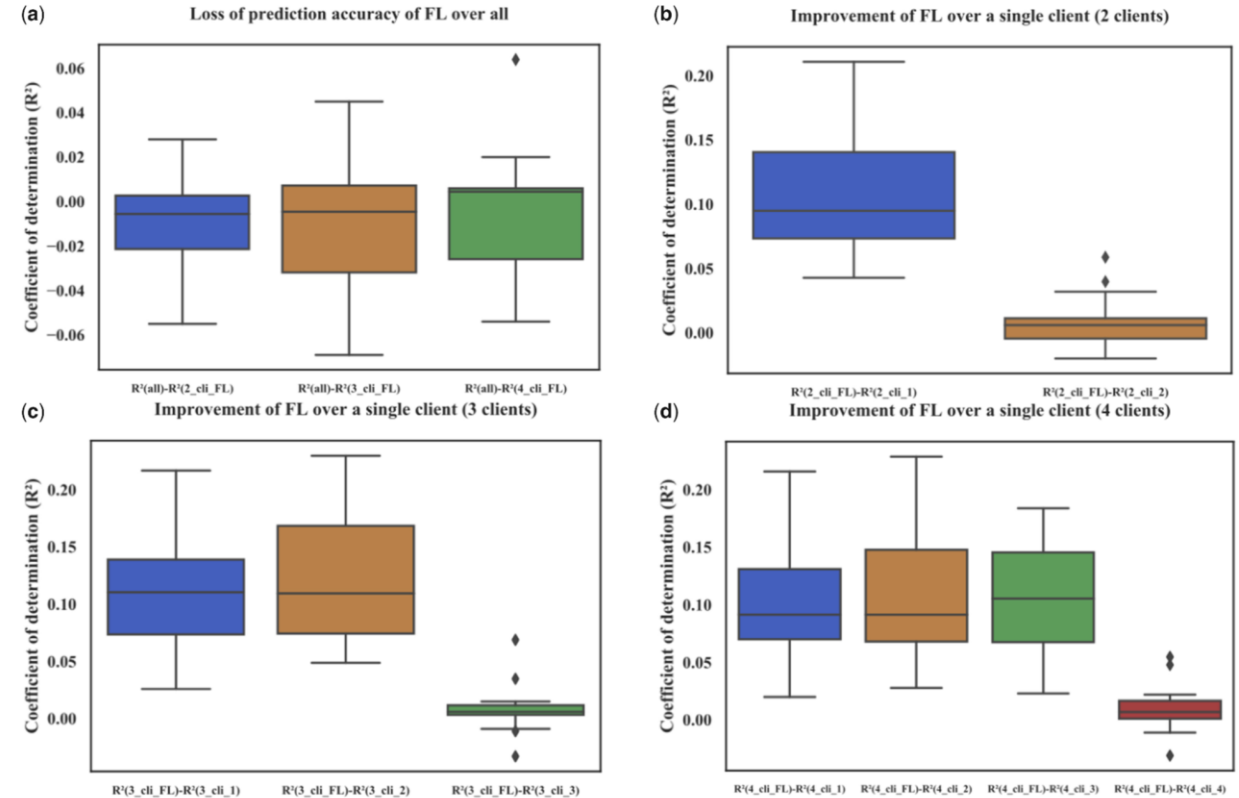
Shaoqi Chen et al. Qiang Yang, Qi Liu, 2020, *Bioinformatics*.

FL-QSAR: Federated Learning Results

Balanced Cases



Imbalanced Cases



Shaoqi Chen et al. Qiang Yang, Qi Liu, 2020, *Bioinformatics*.

Conclusion

- ◆ HFL can achieve almost the same performance as collaboration via cleartext learning algorithms using all shared information.
- ◆ Collaboration by HFL gained a substantially performance improvement than that of one client by using only its private data.
- ◆ A prototype collaborative QSAR system based on FL is presented:
<https://github.com/bm2-lab/FL-QSAR>.

Shaoqi Chen et al. Qiang Yang, Qi Liu, 2020, *Bioinformatics*.

Challenges for Federated Learning

Models [BVH+18]

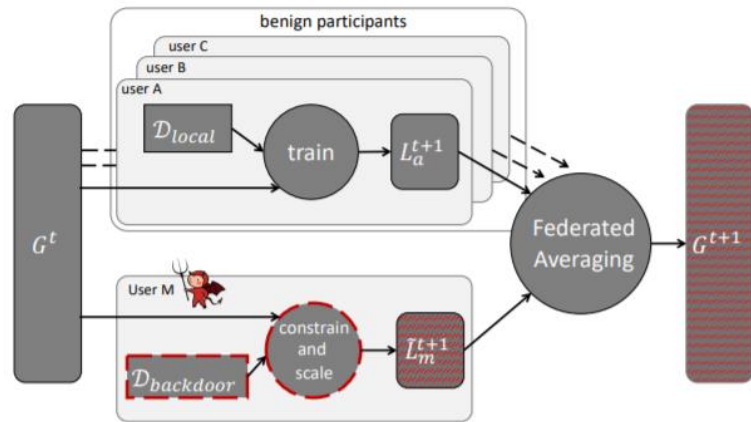
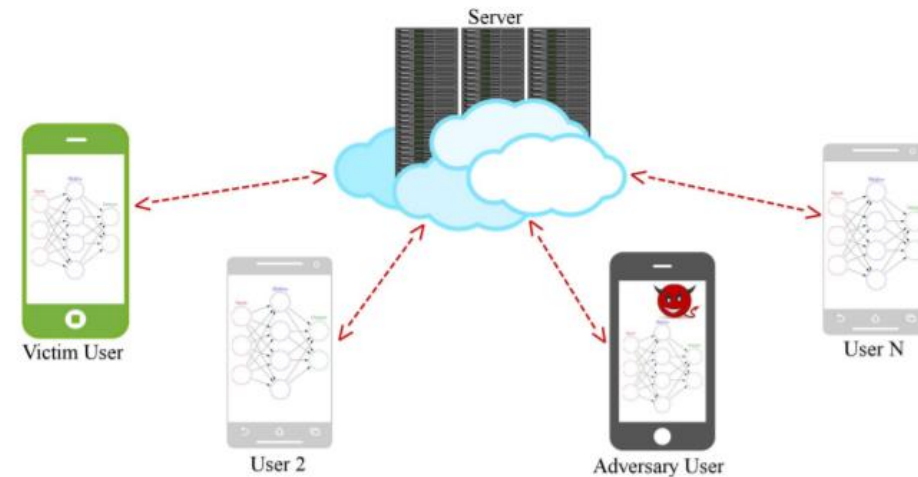


Fig. 1: **Overview of the attack.** The attacker compromises one or more of the participants, trains a model on the backdoor data using our new constrain-and-scale technique, and submits the resulting model. After federated averaging, the global model is replaced by the attacker's backdoored model.

Eugene B et al. 2018. *How To Backdoor Federated Learning*. arXiv:cs.CR/1807.00459

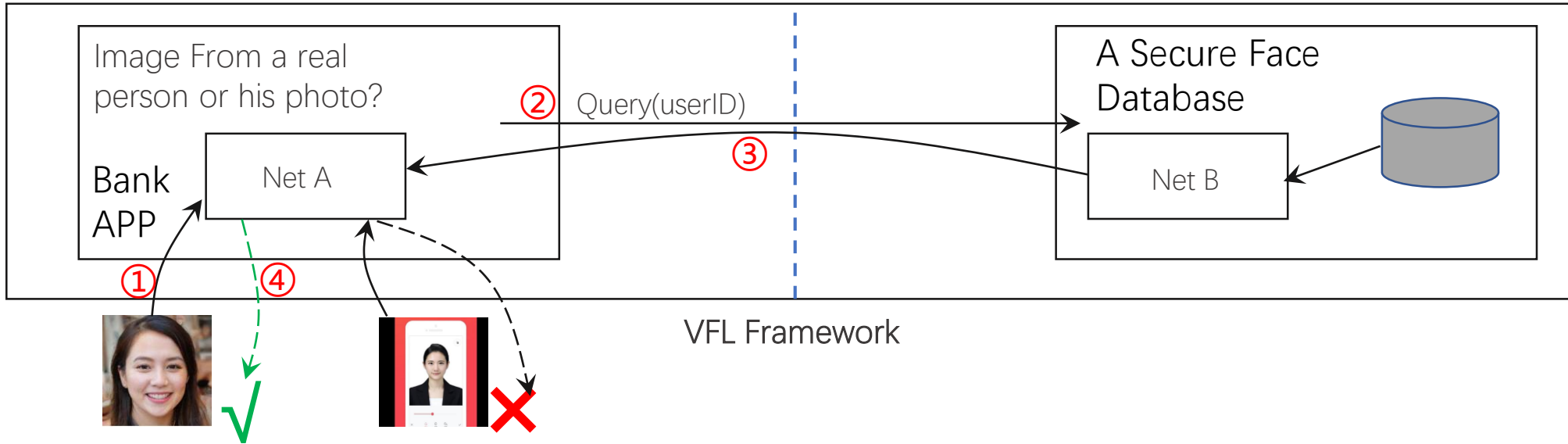
Data [HAP17]



(b) Collaborative Learning

Briland H et al. 2017. *Deep Models Under the GAN: Information Leakage from Collaborative Deep Learning*

New Direction: Auto Vertical Federated Learning



Banking Authentication:

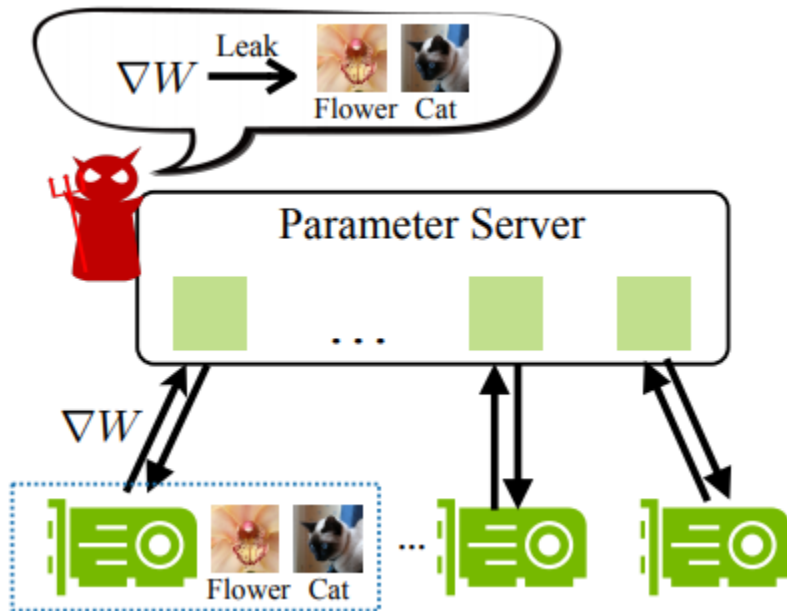
1. Upload front camera photo: to judge whether an image is taken from a real person or his photo.
2. Bank A can cooperative with a **a Secure Face Database** with VFL;

What can AutoFL (Vertical) do:

To determine the learning architecture **automatically** and **locally** in a **communication-efficient manner** with **data protection**;

Privacy Attack Example: Deep Leakage.

Song Han from MIT designed Deep Leakage Attacks that tackle DP-protected models, and are able to reconstruct training data from gradients with pixel-level accuracy.



	Original	$G-10^{-4}$	$G-10^{-3}$	$G-10^{-2}$	$G-10^{-1}$
Accuracy	76.3%	75.6%	73.3%	45.3%	$\leq 1\%$
Defendability	-	✗	✗	✓	✓
		$L-10^{-4}$	$L-10^{-3}$	$L-10^{-2}$	$L-10^{-1}$
Accuracy	-	75.6%	73.4%	46.2%	$\leq 1\%$
Defendability	-	✗	✗	✓	✓

Reconstruct training data



Ground Truth



References

- [McMahan'16] H. Brendan McMahan et al., "Communication-Efficient Learning of Deep Networks from Decentralized Data", Google, 2016.
- [Bonawitz'19] Keith Bonawitz et al., "Towards Federated Learning at Scale: System Design", Google, 2019.
- [Su'18] H. Su and H. Chen. "Experiments on Parallel Training of Deep Neural Network using Model Averaging," Jul. 2018. Available: <https://arxiv.org/abs/1507.01239>
- [Hegedüs'19] I. Hegedüs, G. Danner, and M. Jelasity. "Gossip Learning as a Decentralized Alternative to Federated Learning," In Proceedings of the 14th International Federated Conference on Distributed Computing Techniques, pp. 74–90, Jun. 2019.
- [Daily'18] J. Daily, A. Vishnu, et al., "GossipGrad: Scalable Deep Learning using Gossip Communication based Asynchronous Gradient Descent," Mar. 2018. Available: <http://arxiv.org/abs/1803.05880>
- [Liu'18] Y. Liu, J. Liu and T. Basar. "Differentially Private Gossip Gradient Descent," In IEEE Conference on Decision and Control (CDC'18), pp. 2777–2782, Dec. 2018.
- [Guha'19] N. Guha, A. Talwalkar, and V. Smith. "One-Shot Federated Learning," Mar. 2019. Available: <https://arxiv.org/abs/1902.11175>
- [Phong'19] L.T. Phong, and T.T. Phuong. "Privacy-Preserving Deep Learning via Weight Transmission," Apr. 2019. Available: <https://arxiv.org/abs/1809.03272>
- [Tang'19] H. Tang, C. Yu, C. Renggli, and S. Kassing, et al., "Distributed Learning over Unreliable Networks," May 2019. Available: <https://arxiv.org/abs/1810.07766>
- [Li'19] T. Li, A.K. Sahu, and M. Zaheer, et al., "Federated Optimization for Heterogeneous Networks," Jul. 2019. <https://arxiv.org/abs/1812.06127>
- [Chen'19] X. Chen, T. Chen, and H. Sun, et al., "Distributed training with heterogeneous data: Bridging median- and mean-based algorithms," Jun. 2019. Available: <https://arxiv.org/abs/1906.01736>
- [Liu'18] D. Liu, T. Miller, et al., "FADL: Federated-Autonomous Deep Learning for Distributed Electronic Health Record," Dec. 2018. Available: <https://arxiv.org/abs/1811.11400>
- [Bagdasaryan'19] E. Bagdasaryan, A. Veit, and Y. Hua, et al., "How to backdoor federated learning," arXiv preprint arXiv:1807.00459, Aug. 2019. Available: <https://arxiv.org/abs/1807.00459>
- [Bhagoji'19] A. N. Bhagoji, S. Chakraborty, et al., "Analyzing federated learning through an adversarial lens," arXiv preprint arXiv:1811.12470, Mar. 2019. Available: <http://arxiv.org/abs/1811.12470>
- [Lin'18] Y. Lin, S. Han, et al., "Deep Gradient Compression: Reducing the Communication Bandwidth for Distributed Training," Feb. 2018. Available: <https://arxiv.org/abs/1712.01887>
- [Reisizadeh '19] A. Reisizadeh, A. Mokhtari, et al., "FedPAQ: A Communication-Efficient Federated Learning Method with Periodic Averaging and Quantization," Oct. 2019. Available: <https://arxiv.org/abs/1909.13014>
- [Wang'19] L. Wang, W. Wang, and B. Li, "CMFL: Mitigating Communication Overhead for Federated Learning," Jul. 2019. Available: <https://www.cse.ust.hk/~weiwa/papers/cmfl-icdcs19.pdf>
- [Konečný'17], H. Brendan McMahan, et al., "Federated Learning: Strategies for Improving Communication Efficiency," Oct. 2017. Available: <https://arxiv.org/abs/1903.0742>
- [Bonawitz'16] K. Bonawitz, et al., "Practical Secure Aggregation for Federated Learning on User-Held Data," Nov. 2016. Available: <https://arxiv.org/abs/1611.0448>
- [Kang'19] J. Kang, Z. Xiong, et al., "Incentive Design for Efficient Federated Learning in Mobile Networks: A Contract Theory Approach," May 2019. Available: <https://arxiv.org/abs/1905.07479>
- [Richardson'19] A. Richardson, A. Filos-Ratsikas, and B. Faltings, "Rewarding High-Quality Data via Influence Functions," Aug. 2019. Available: <https://arxiv.org/abs/1908.11598>
- [Feng'18] S. Feng, D. Niyato, et al., "Joint Service Pricing and Cooperative Relay Communication for Federated Learning," Nov. 2018. Available: <https://arxiv.org/abs/1811.12082>
- [Wang'19] G. Wang, C. Dang, and Z. Zhou, "Measure Contribution of Participants in Federated Learning," Sep. 2019. <https://arxiv.org/abs/1909.0852>