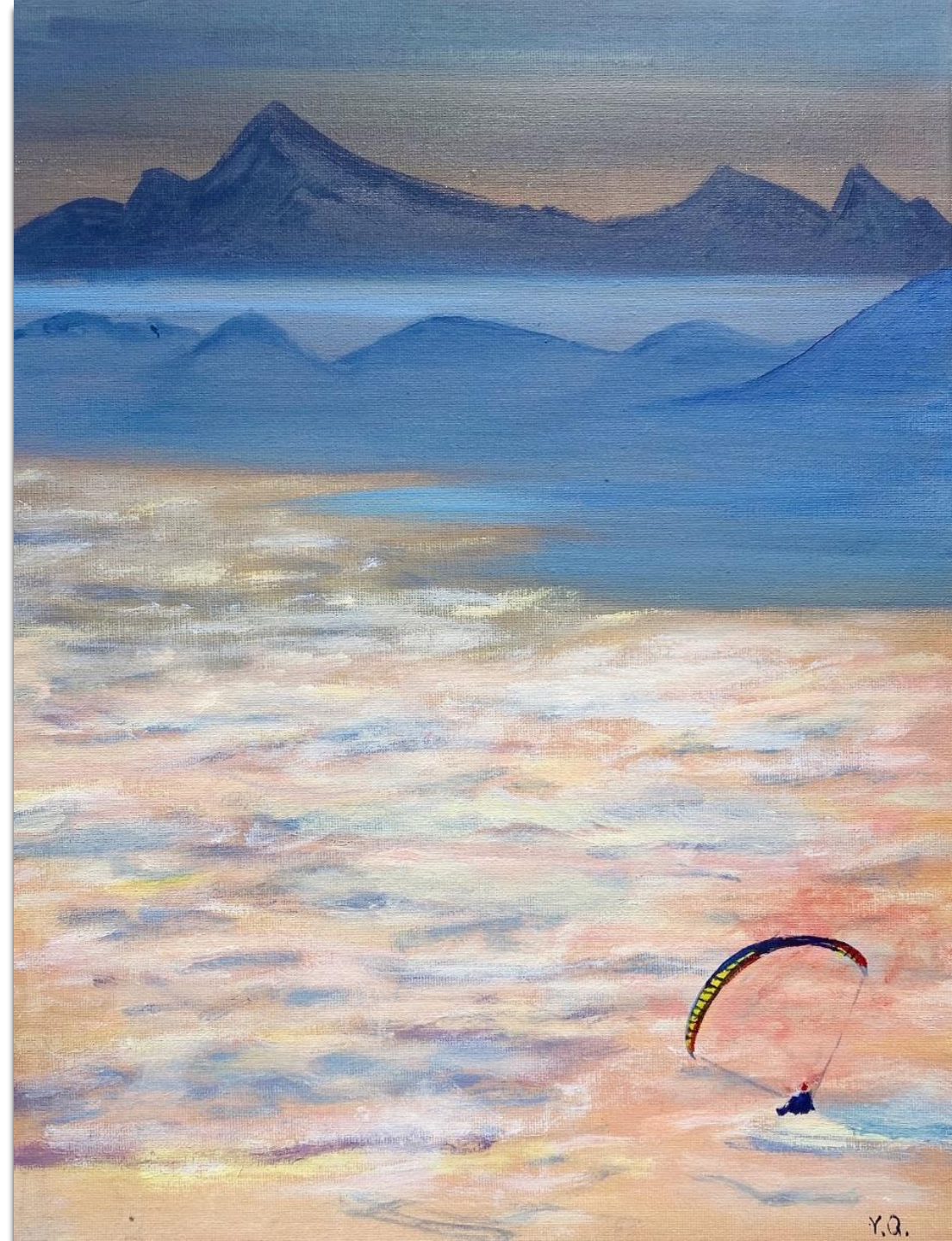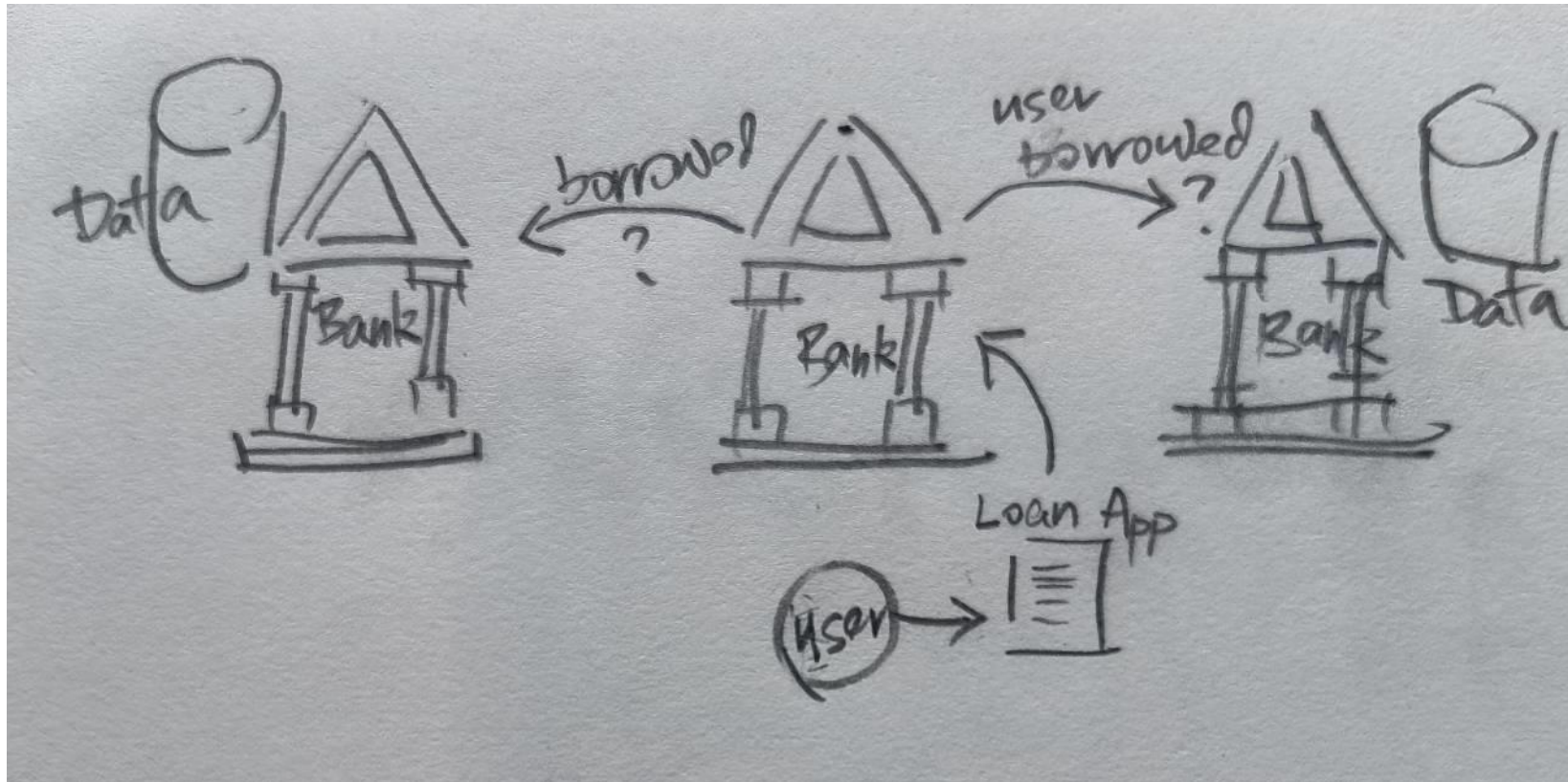# Trustworthy Federated Learning
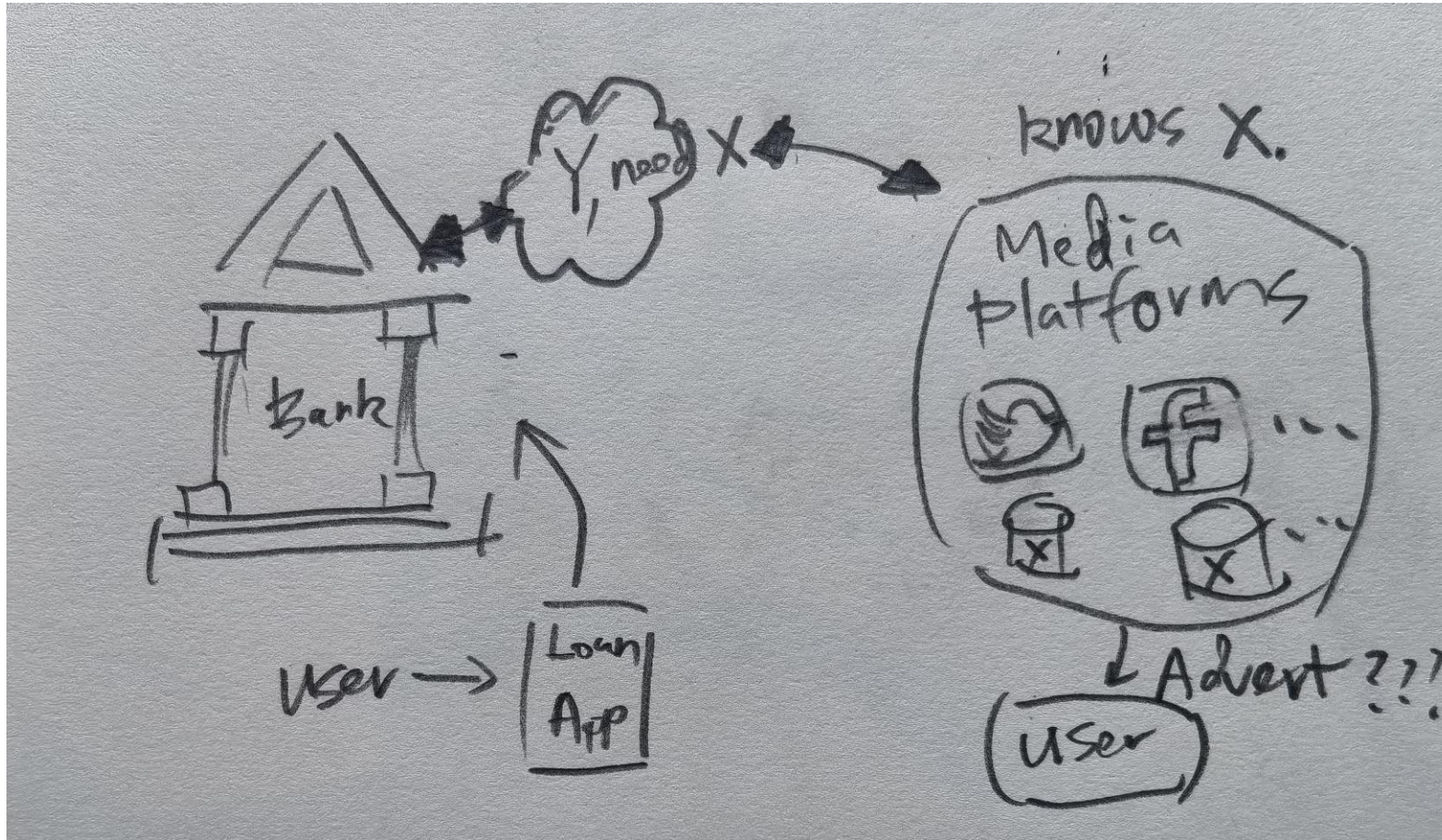
**QIANG YANG**

HKUST，WeBank

WeBank

# Case 1: Bank Loan, Anti-fraud Applications



- User Applying Bank Loan
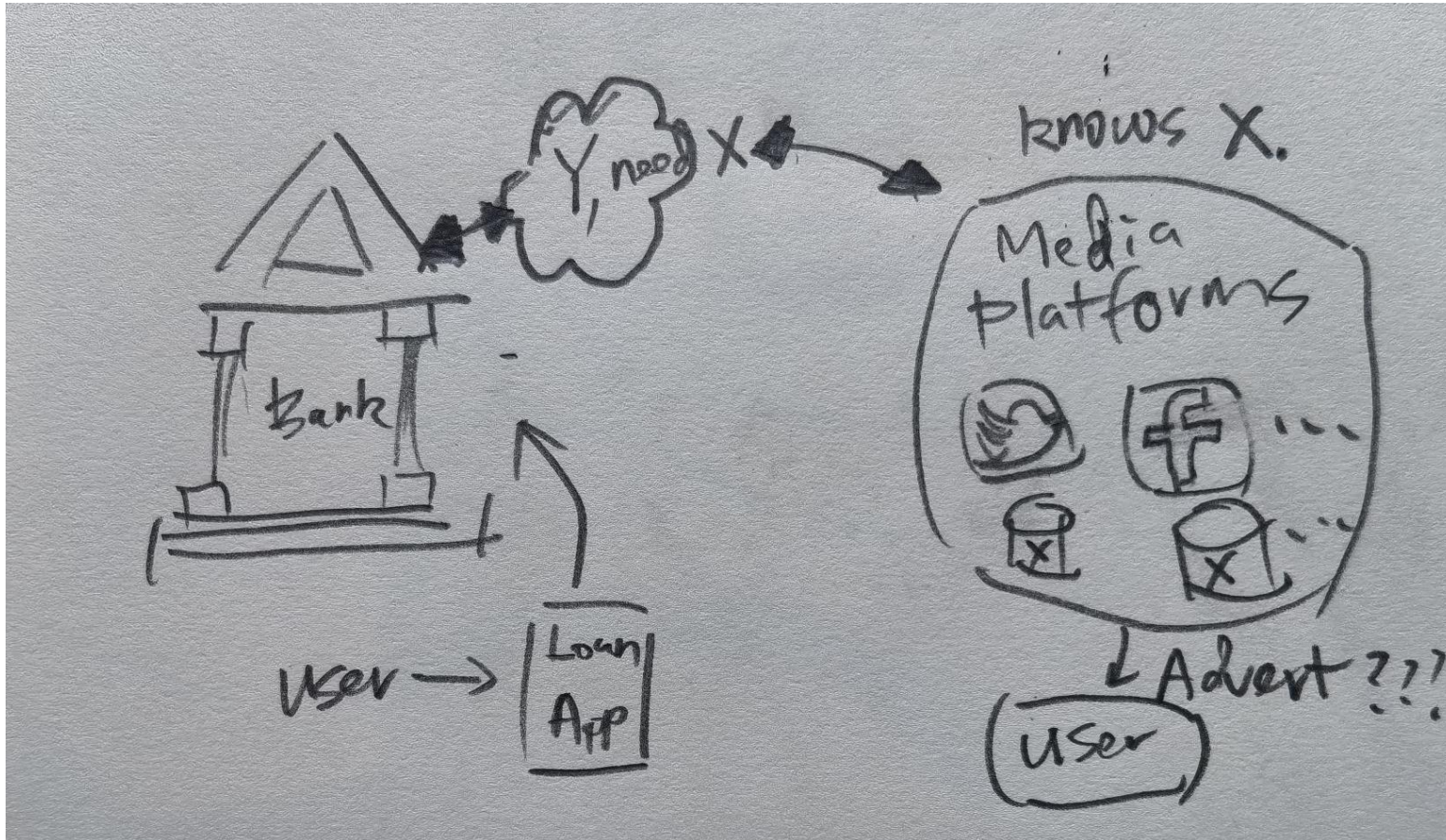  - Has the user got loans from other banks?
  - Is this a case of fraud?

# Case 2: Financial Credit Rating



- User Applying Bank Loan
  - Is the user credible?
  - How is the user's other behavior?
- Bank knows Y (credit history)
- Media knows X (Behavior)
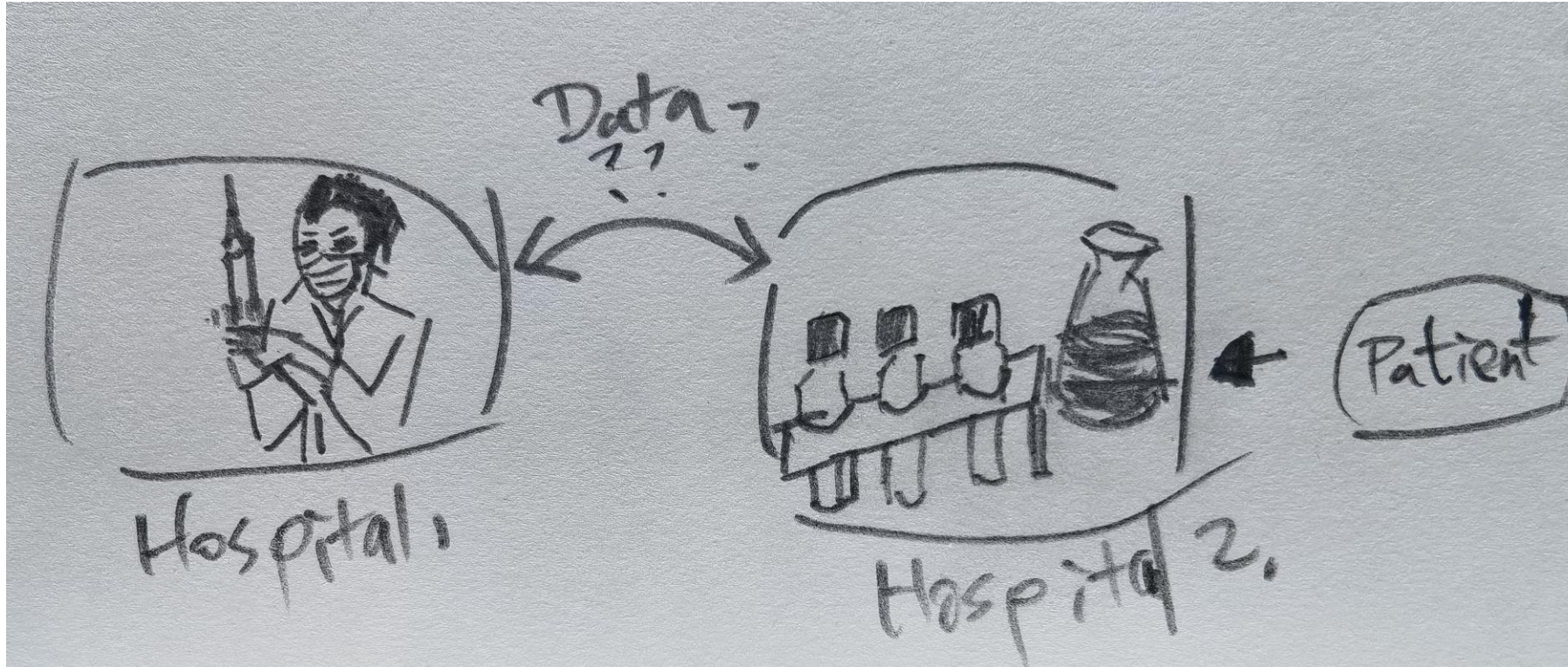- Can they build a collective model at the bank side?

# Case 3: Advertising



- Bank needs to advertise to users with good credit history
- But the data cannot be given to media companies
- Should a media company advertise to a user?

WeBank

# Case 4: Multiple hospitals build collective model



- Early detection of illness such as stroke
- Design new drugs
- ...

Qiang Yang 202212

# 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

- Three major data security laws in 2021:
  - Personal Data Privacy Protection Regulation
  - Data Security Law
  - Cybersecurity Law
- Require businesses to get user consent, not leak out or tamper with data
- When conducting data transactions with third parties, proposed contract should follow legal data protection regulations.

## 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:

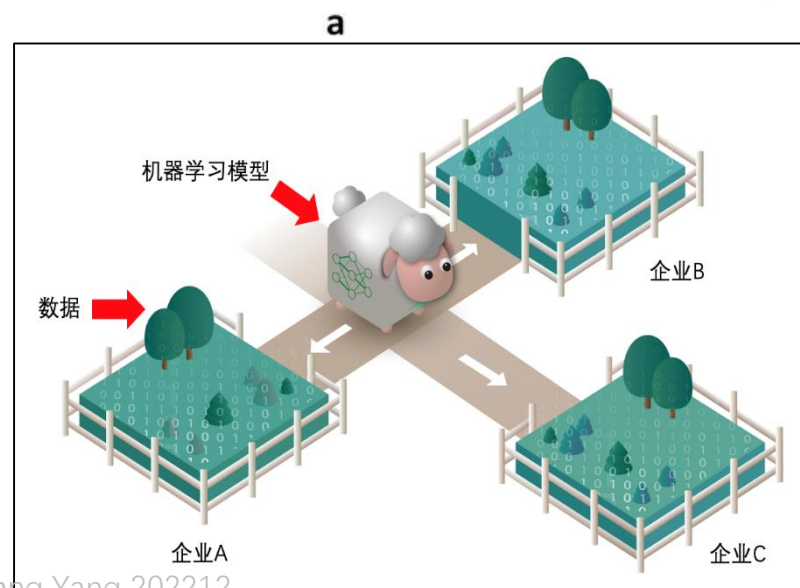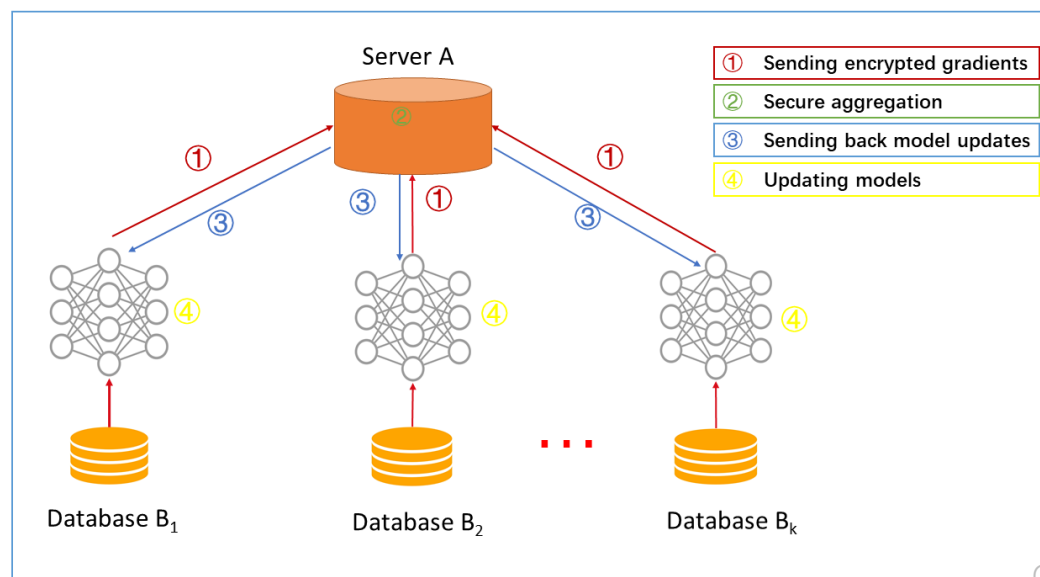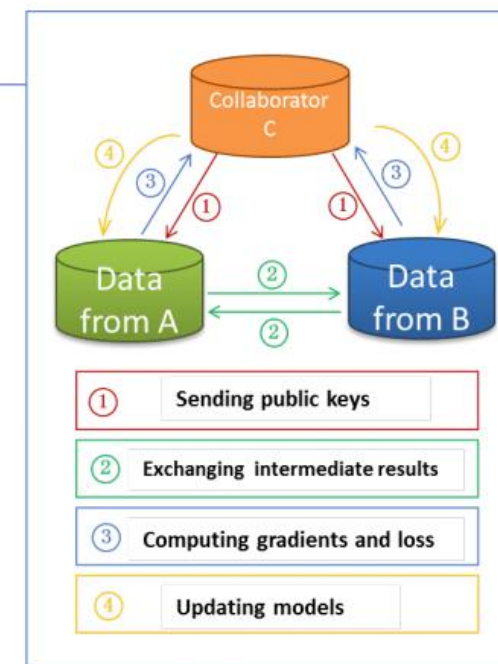| | |
|---|---|
| **Personal information protection** | The Cybersecurity Law clearly states requirements for the collection, use and protection of personal information. |
| **Critical information infrastructure** | The Cybersecurity Law frequently mentions the protection of "critical information infrastructure". |
| **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. |

**From Report by KPMG 2017**

WeBank

# **Federated Learning**

- Federated Machine Learning
  - Multiple parties compute a model
  - Each party contributes some local data, without sharing the data
- Horizontal and Vertical Federated Learning

# Key Components in Federated Learning

Model Design

Communication

Privacy &
Security

Distributed ML

Incentive

Image from Google

# Horizontal Federated Learning: Divide by Users/Samples



(a) Horizontal Federated Learning

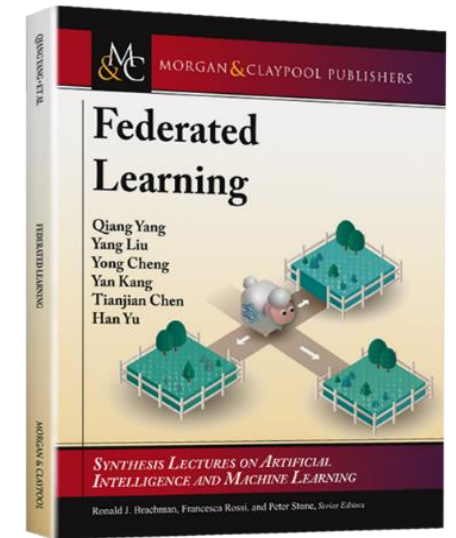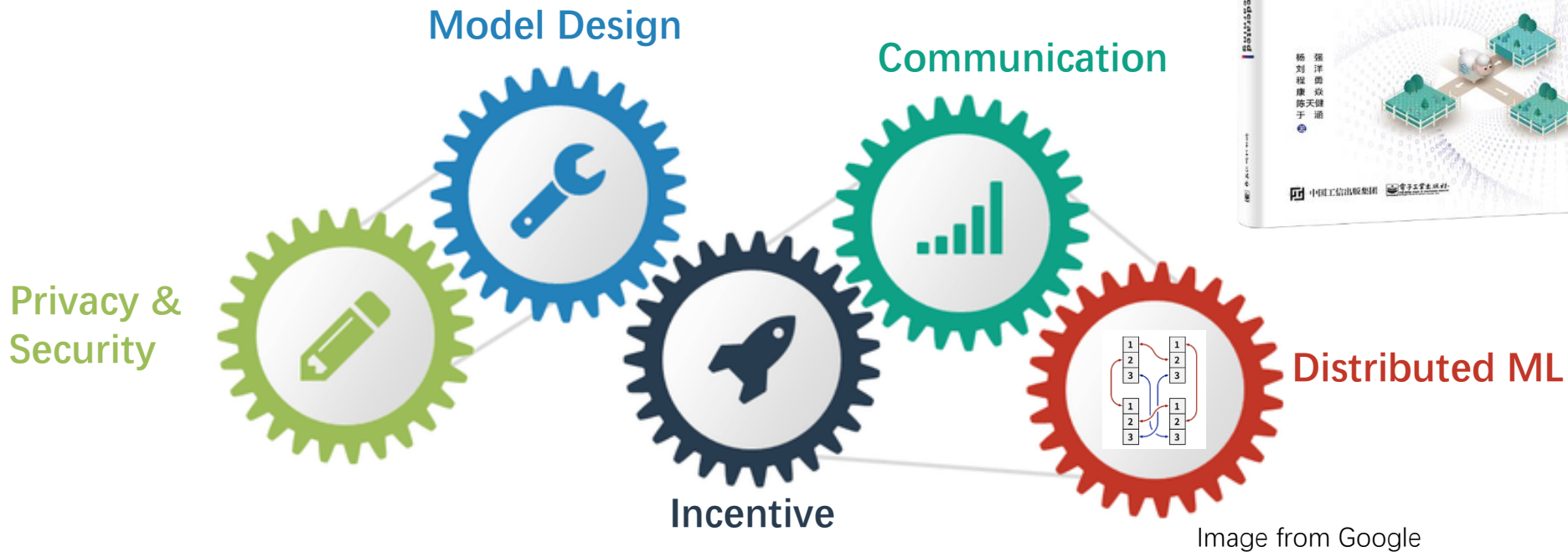**FEDERATED LEARNING FOR MOBILE KEYBOARD PREDICTION, Andrew Hard, et al., Google, 2018**

**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

# 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).
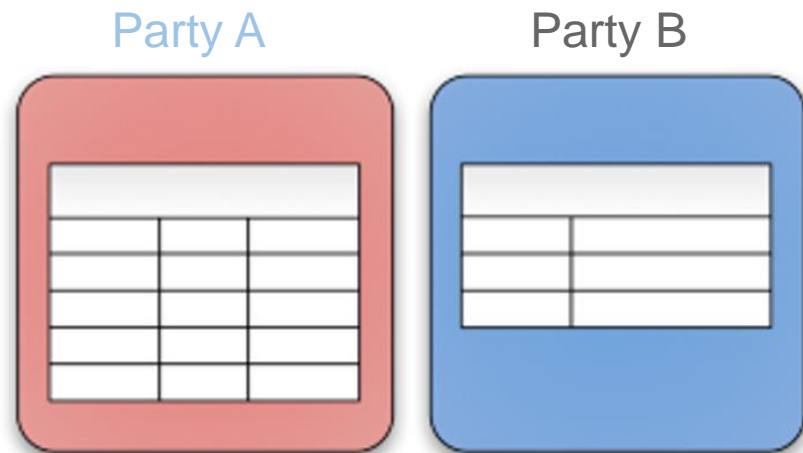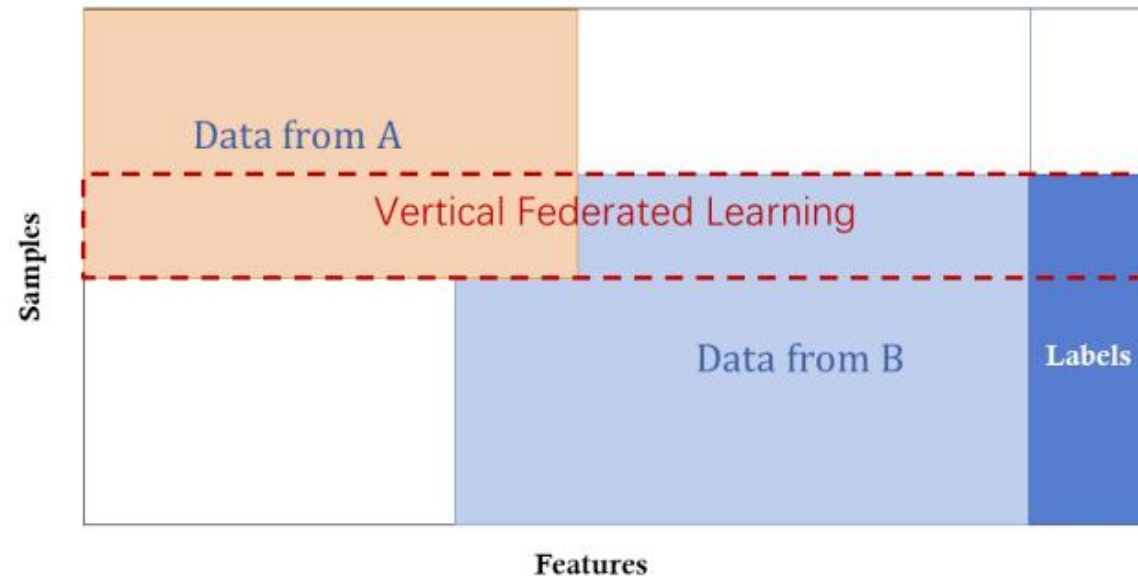
Party A          Party B



Image from Google

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.

# XGBoost in Federated Learning



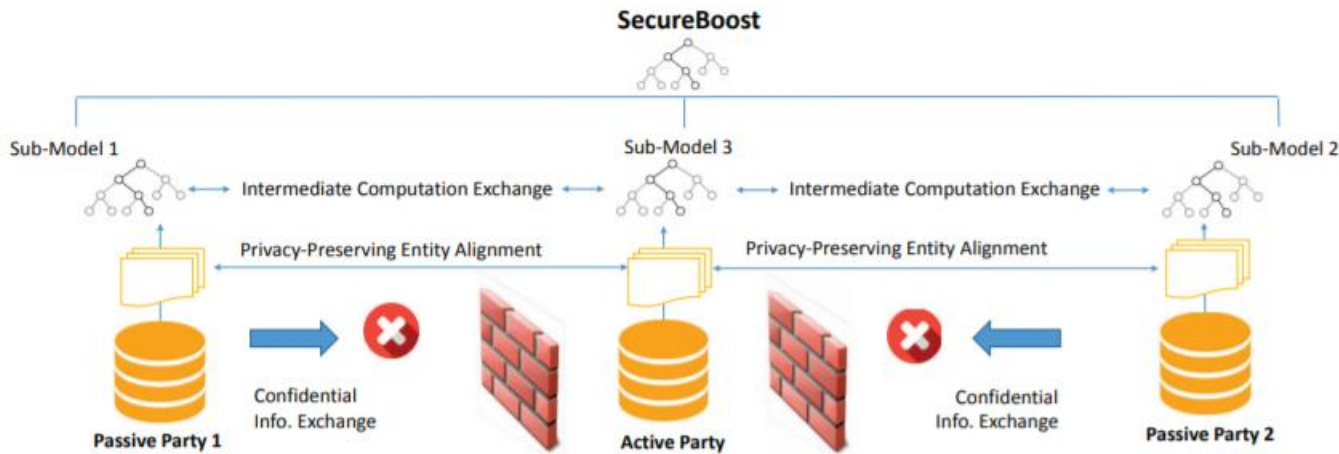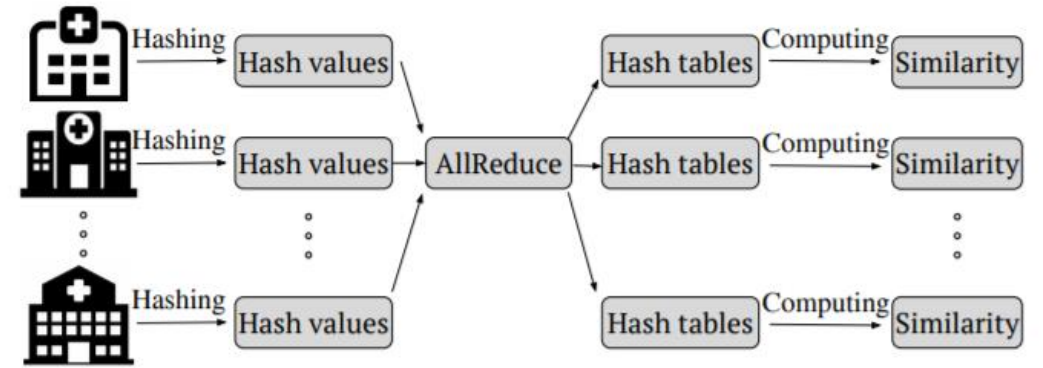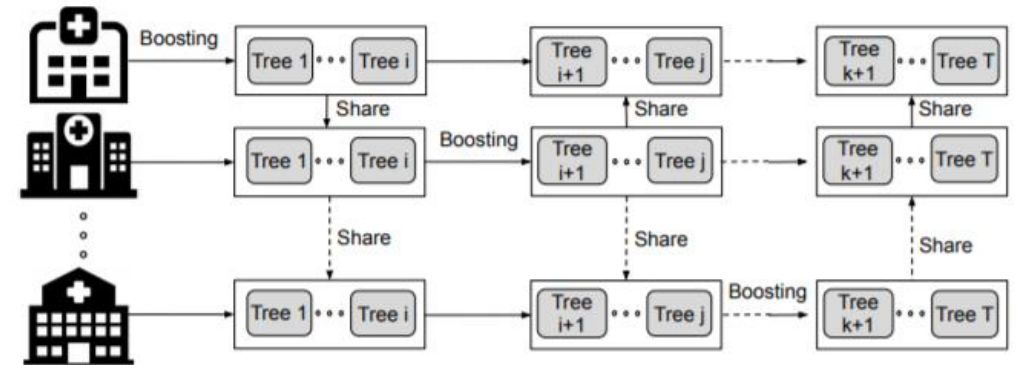Figure 1: Illustration of the proposed SecureBoost framework

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

WeBank

# 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?



Data Owner 1

Dataset

FL Model

Revenue from FL customers

Limited Budget for Utility Transfer to data owners at $t$, $B(t)$

Data Owner $i$ ... ... Data Owner $N$

$u_1(t)$  $u_N(t)$  $u_i(t)$

•Qiang Yang, Yang Liu, Tianjian Chen, Yongxin Tong:
Federated Machine Learning: Concept and Applications. ACM TIST 10(2): 12:1-12:19 (2019)

# Federated Learning: Best for AI Models

## One-off MPC Projects



## Standard Federated Learning

encryption

**Federated Learning**

Trusted execution environment (hardware)

Secret Sharing

Differential privacy

**WeBank**

# Advances and Open Problems in Federated Learning

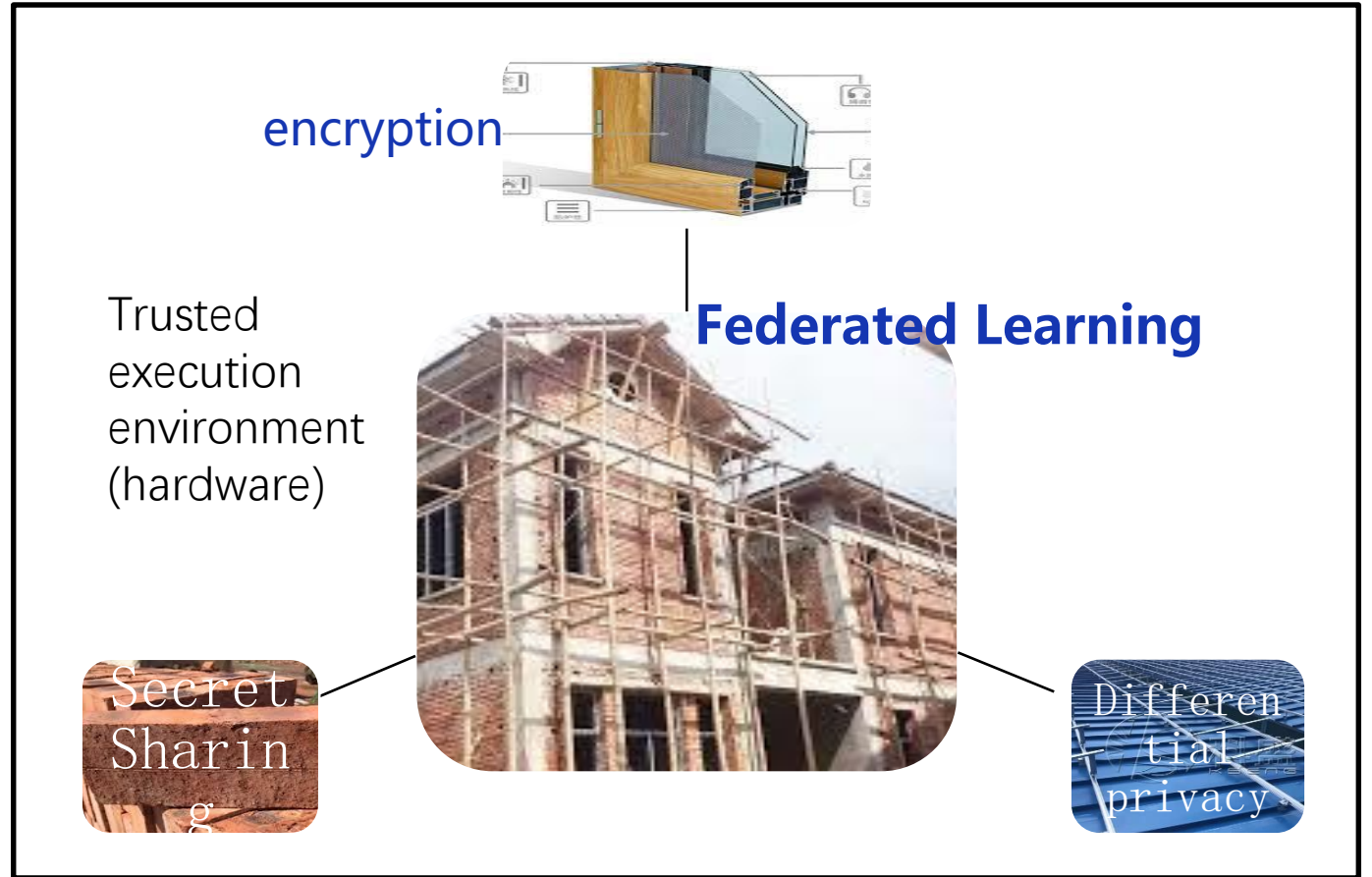Peter Kairouz[7*]     H. Brendan McMahan[7*]     Brendan Avent[21]     Aurélien Bellet[9]

Mehdi Bennis[19]     Arjun Nitin Bhagoji[13]     Keith Bonawitz[7]     Zachary Charles[7]

Graham Cormode[23]     Rachel Cummings[6]     Rafael G.L. D'Oliveira[14]

Salim El Rouayheb[14]     David Evans[22]     Josh Gardner[24]     Zachary Garrett[7]

Adrià Gascón[7]     Badih Ghazi[7]     Phillip B. Gibbons[2]     Marco Gruteser[7,14]

Zaid Harchaoui[24]     Chaoyang He[21]     Lie He [4]     Zhouyuan Huo [20]

Ben Hutchinson[7]     Justin Hsu[25]     Martin Jaggi[4]     Tara Javidi[17]     Gauri Joshi[2]

Mikhail Khodak[2]     Jakub Konečný[7]     Aleksandra Korolova[21]     Farinaz Koushanfar[17]

Sanmi Koyejo[7,18]     Tancrède Lepoint[7]     Yang Liu[12]     Prateek Mittal[13]

Mehryar Mohri[7]     Richard Nock[1]     Ayfer Özgür[15]     Rasmus Pagh[7,10]

Mariana Raykova[7]     Hang Qi[7]     Daniel Ramage[7]     Ramesh Raskar[11]

Dawn Song[16]     Weikang Song[7]     Sebastian U. Stich[4]     Ziteng Sun[3]

Ananda Theertha Suresh[7]     Florian Tramèr[15]     Praneeth Vepakomma[11]     Jianyu Wang[2]

Li Xiong[5]     Zheng Xu[7]     Qiang Yang[8]     Felix X. Yu[7]     Han Yu[12]     Sen Zhao[7]

[1]Australian National University, [2]Carnegie Mellon University, [3]Cornell University,
[4]École Polytechnique Fédérale de Lausanne, [5]Emory University, [6]Georgia Institute of Technology,
[7]Google Research, [8]Hong Kong University of Science and Technology, [9]INRIA, [10]IT University of Copenhagen,
[11]Massachusetts Institute of Technology, [12]Nanyang Technological University, [13]Princeton University,
[14]Rutgers University, [15]Stanford University, [16]University of California Berkeley,
[17] University of California San Diego, [18]University of Illinois Urbana-Champaign, [19]University of Oulu,
[20]University of Pittsburgh, [21]University of Southern California, [22]University of Virginia,
[23]University of Warwick, [24]University of Washington, [25]University of Wisconsin–Madison
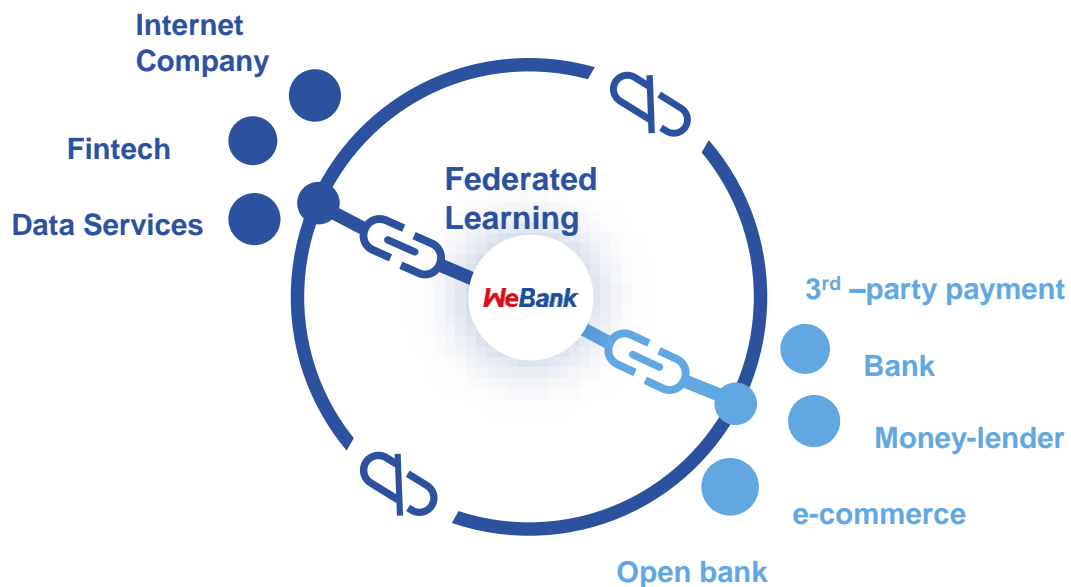
# Multiple Data Sources in Credit Rating

**Build local models** ---> **Federated Modeling** ---> **Rating**

Internet Company

Fintech

Data Services

Federated Learning

**WeBank**

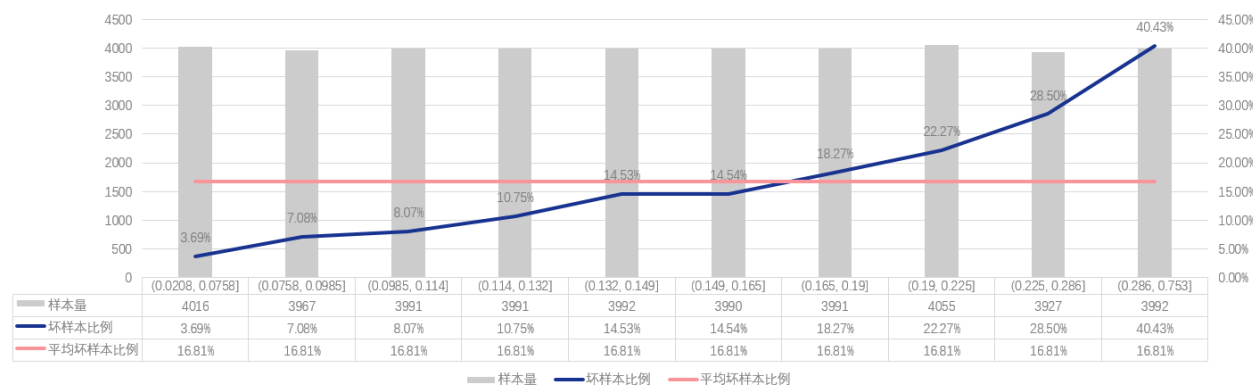3rd –party payment

Bank

Money-lender

e-commerce

Open bank

**Increase the accuracy of models**

**Score to Indicate High/Medium/Low Risk Level**

Model performance :
   AUC=0.70
   KS=30
   x2.4 bad rate in the bottom decile

| | (0.0208, 0.0758] | (0.0758, 0.0985] | (0.0985, 0.114] | (0.114, 0.132] | (0.132, 0.149] | (0.149, 0.165] | (0.165, 0.19] | (0.19, 0.225] | (0.225, 0.286] | (0.286, 0.753] |
|---|---|---|---|---|---|---|---|---|---|---|
| 样本量 | 4016 | 3967 | 3991 | 3991 | 3992 | 3990 | 3991 | 4055 | 3927 | 3992 |
| 坏样本比例 | 3.69% | 7.08% | 8.07% | 10.75% | 14.53% | 14.54% | 18.27% | 22.27% | 28.50% | 40.43% |
| 平均坏样本比例 | 16.81% | 16.81% | 16.81% | 16.81% | 16.81% | 16.81% | 16.81% | 16.81% | 16.81% | 16.81% |

样本量   坏样本比例   平均坏样本比例

WeBank

# Construction-Site Safety w/ Federated Computer Vision

## WeBank AI X Extreme Vision

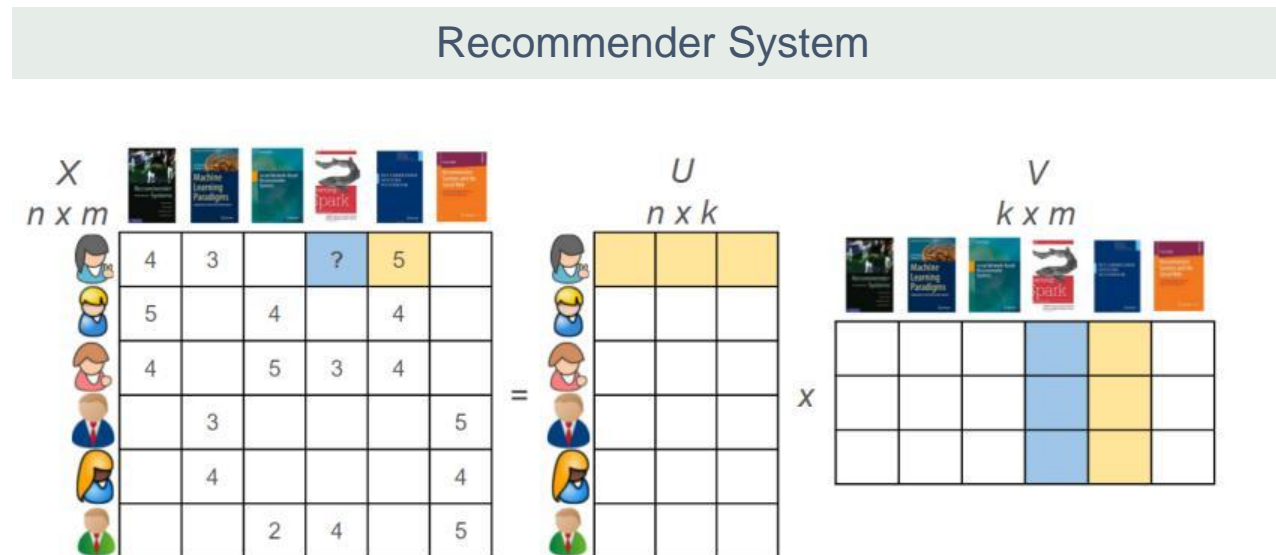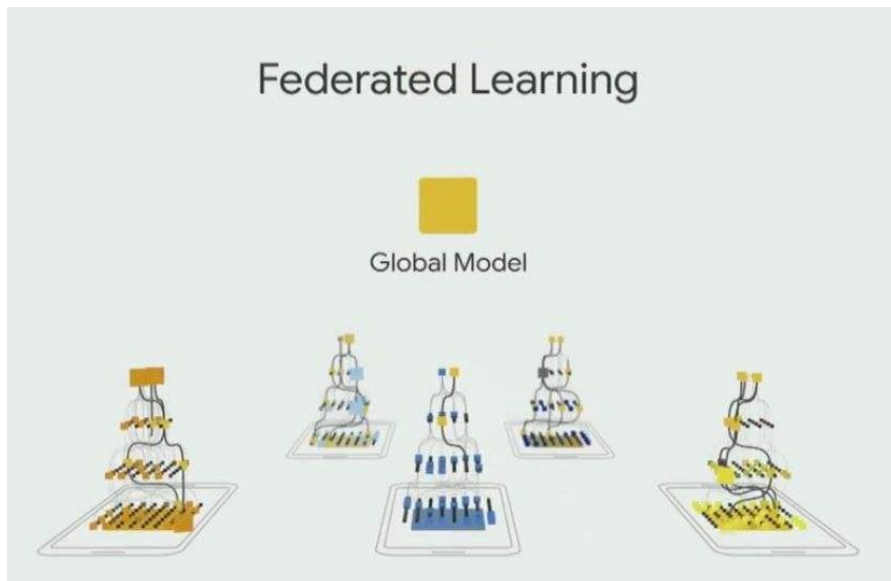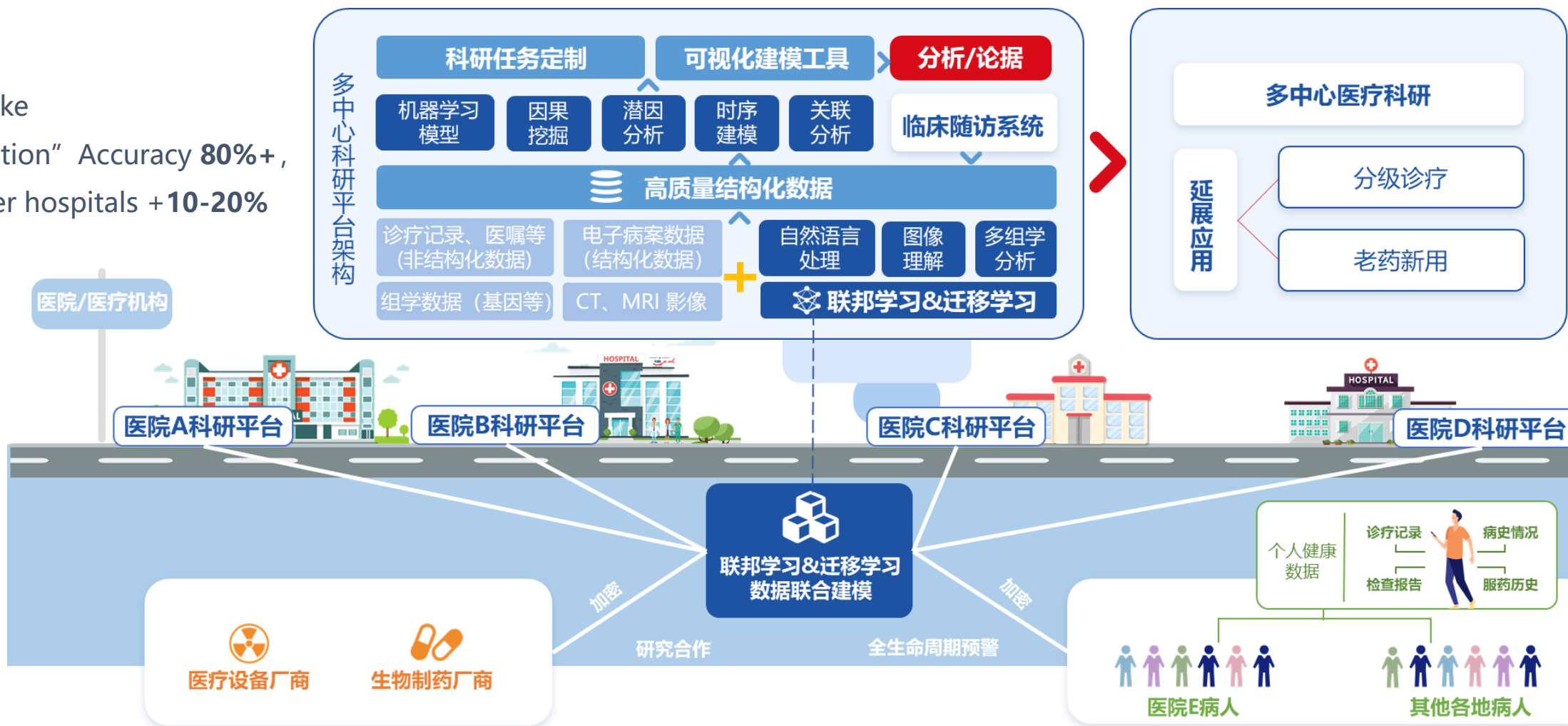# Federated Recommendation Systems



**Assumption**: a trustworthy 3rd-party as coordinator, which can be removed
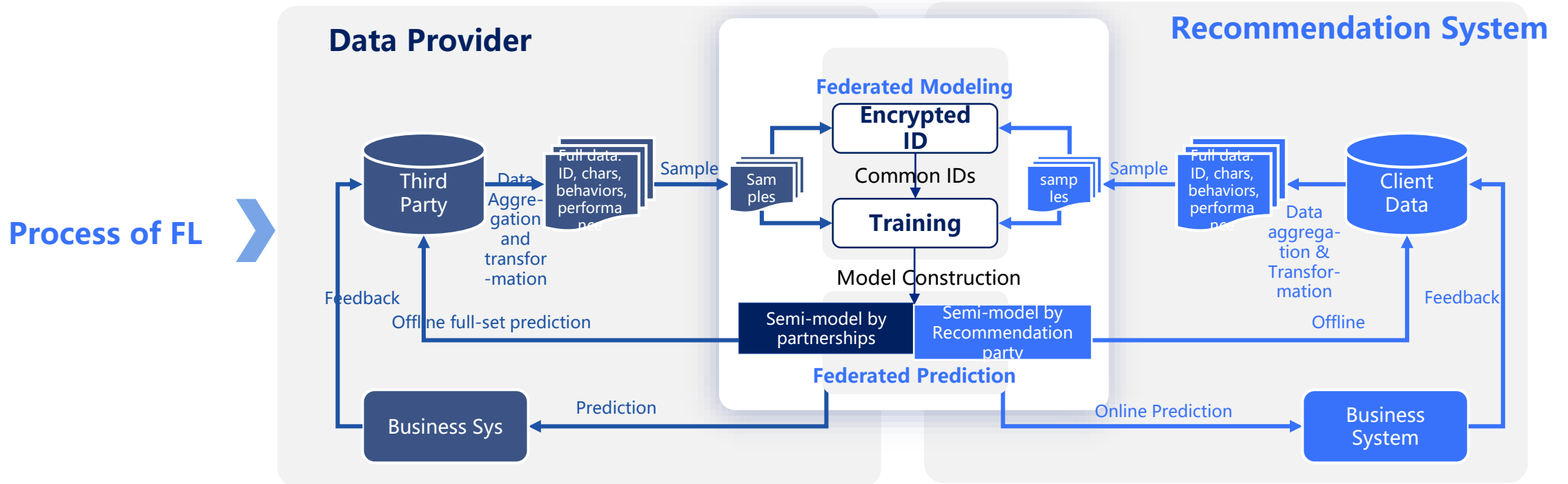
# Drug Discovery from Multiple Data Sources

Multiple hospitals, medical centers and drug companies

- "Stroke Prediction" Accuracy **80%+** , Smaller hospitals +**10-20%**

# News Recommendation on Mobile Phones



**Mobile APP** → Multiple Data / Reinforcement Learning → **Operations** → Federated / Recommendation → **Articles** →

Page View + **21%**

Reading Time + **22%**

Click Through +**11%**

**Process of FL**

**Data Provider**

**Recommendation System**

**Federated Modeling**

Third Party

Full data. ID, chars, behaviors, performance — Sample → Samples

**Encrypted ID**

Common IDs

**Training**

samples ← Sample — Full data. ID, chars, behaviors, performance ← Client Data

Data Aggre-gation and transfor-mation

Data aggrega-tion & Transfor-mation

Model Construction

Semi-model by partnerships | Semi-model by Recommendation party

**Federated Prediction**

Feedback

Offline full-set prediction

Offline

Feedback

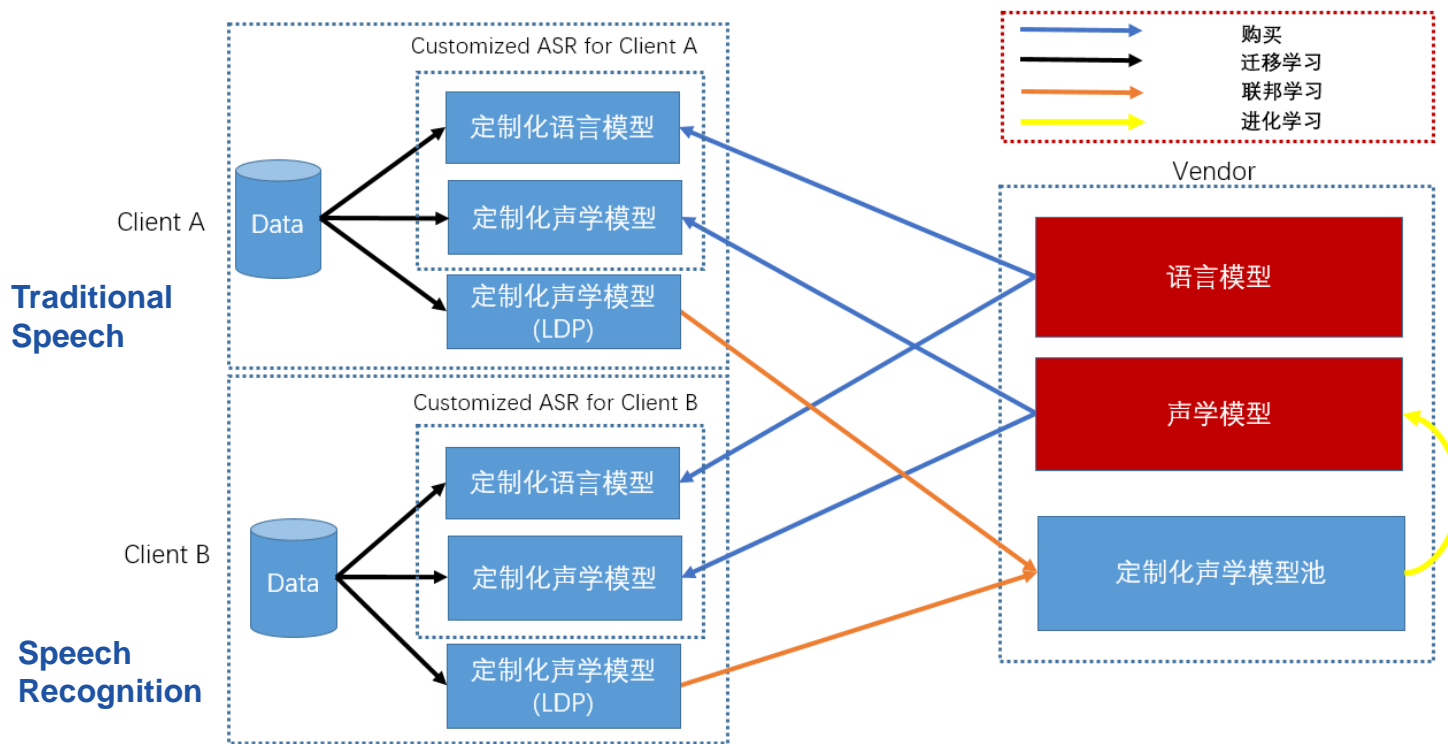Business Sys ← Prediction

Online Prediction → Business System

Note: Customer ID includes but not limited to government ID, mobile phone numbers, device ID, etc.;
Federated model training is initiated by the party who has y-feature

Qiang Yang 202212

**WeBank**

# Speech Recognition on Multiple Data Sources



**1**

**Vendor仅可以根据Client提供的文本信息微调语言模型**

**Speech Models**

**2**

**Client需要给Vendor暴露明文数据来微调ASR**

**Client Privacy Protection**

**3**

**Vendor和Client之间只存在购买这一单向行为**

**Recurrent training**

TFE= Transfer + Federated + Evolutionary

**Traditional Speech**

**Speech Recognition**

Customized ASR for Client A

Client A — Data — 定制化语言模型 / 定制化声学模型 / 定制化声学模型 (LDP)

Customized ASR for Client B

Client B — Data — 定制化语言模型 / 定制化声学模型 / 定制化声学模型 (LDP)

购买
迁移学习
联邦学习
进化学习

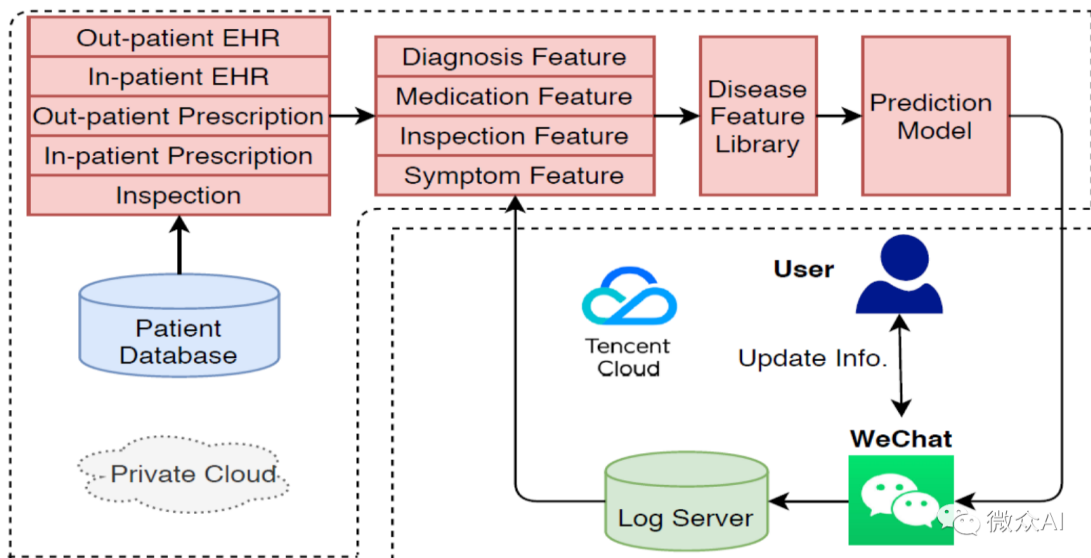Vendor — 语言模型 / 声学模型 / 定制化声学模型池

**Reduced error by 10%~20%**

WeBank

# WeBank-AI and Tencent Joint Project

- Accuracy: Improvement over 80%

- Small hospitals: 10-20%

## FL-IJCAI' 20 Paper

**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 Zhao[2], Yang Liu[1], Hongshan Li[3], Tianjian Chen[1], Xinwei Zhang[4], Dashan Gao[5,6], Ben Tan[1], Han Yu[7] and Yuan Jin[8]

# Application：online service robots KYC, AI services @ WeBank）

## AI soft robots training of service persons

### 智能培训

通过"语音识别+语义理解+语音合成+智能质检"方案，录入真实案例，**以场景互动方式传授说话技巧并与学员进行对话练习。**



**功能优势**

- 录入海量教学案例以语音对话方式进行教学
- 模拟真实销售、客服服务场景与学员情景对练
- 对学员的回复进行实时质检
- 建立学员画像，对学员成绩进行个性化分析
- 支持APP、小程序、PC端多渠道教学

**应用场景&效果**

面对各业务销售及客服服务领域，提供高效、便捷的情景对练式培训，有效降低人力成本，提升培训效果，降低销售、客服人员流动性大带来的培训风险。



保险销售　　汽车销售　　房产销售　　客服培训

## AI soft robots answering questions online

### 电话坐席助手

通过语音识别、语义理解、用户画像等技术，**坐席助手可以提供用户画像辅助坐席决策，提供实时话术推荐辅助坐席对话，并对所有客服通话内容进行实时质检和管控。**



**功能优势**

- 语音内容实时转译
- 实时知识点推荐搜索
- 业务流程提示，标准话术引导
- 对通话内容实时质检
- 用户画像信息展示

**应用场景&效果**

将纯人工电话坐席场景应用模式升级，深度人机协同，实时过程管控，端到端话术沉淀。



业务效率　　服务风险　　营销转化

# Is the Gradient Info Safe to Share?

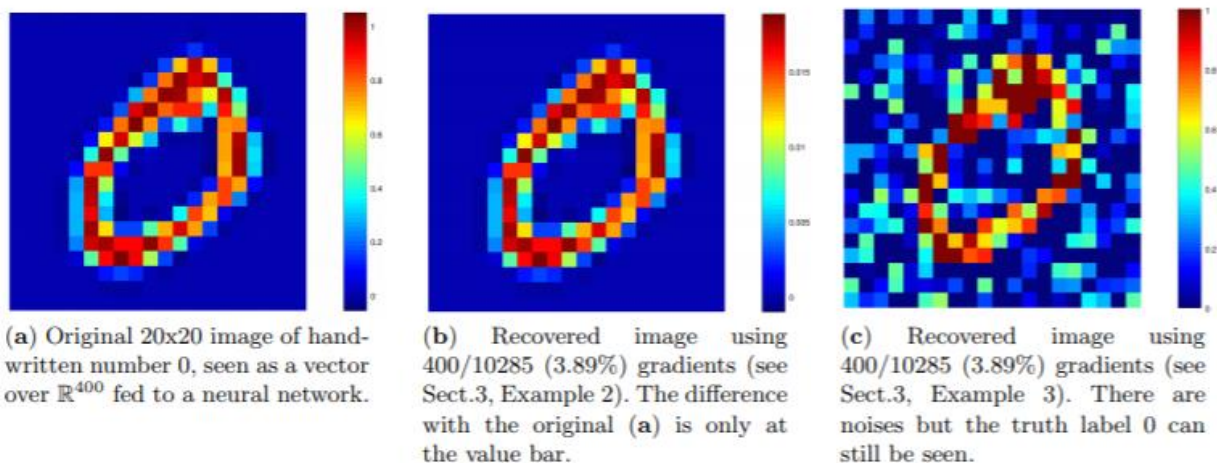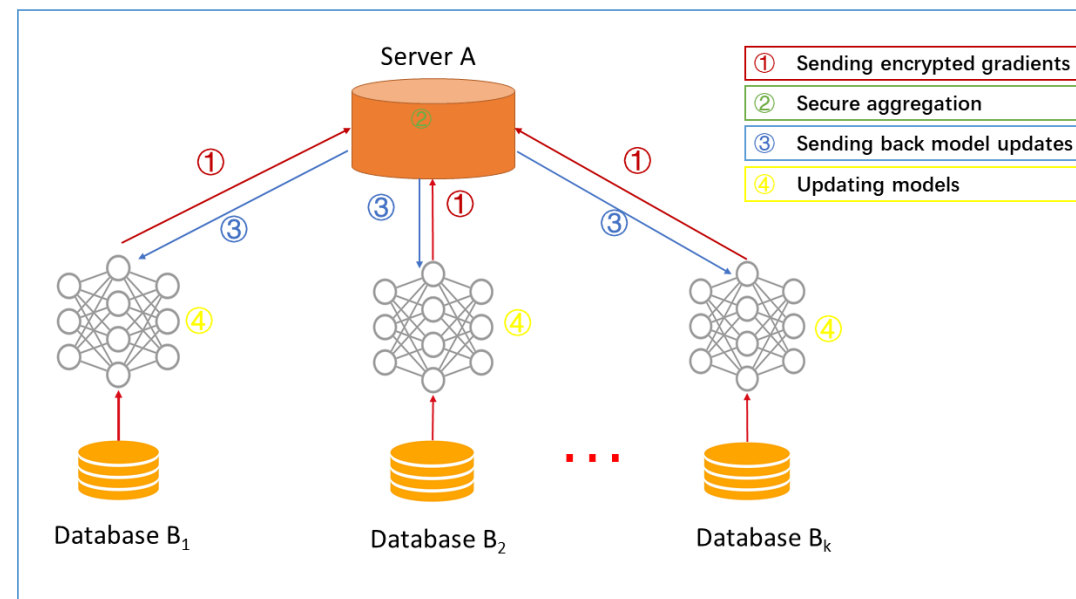*Protect gradients with Homomorphic Encryption*



(a) Original 20x20 image of hand-written number 0, seen as a vector over $\mathbb{R}^{400}$ fed to a neural network.

(b) Recovered image using 400/10285 (3.89%) gradients (see Sect.3, Example 2). The difference with the original (a) is only at the value bar.

(c) Recovered image using 400/10285 (3.89%) gradients (see Sect.3, Example 3). There are noises but the truth label 0 can still be seen.

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



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

**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

**WeBank**

# Challenges for Federated Learning

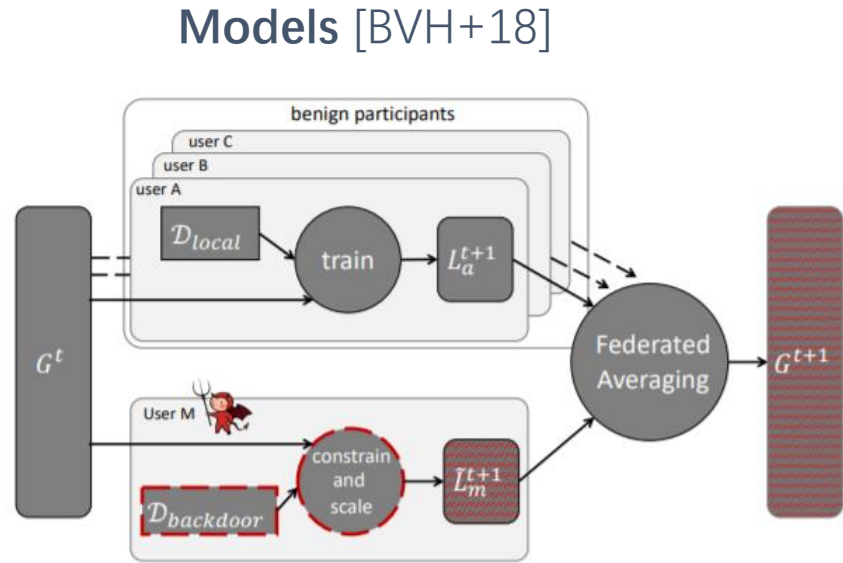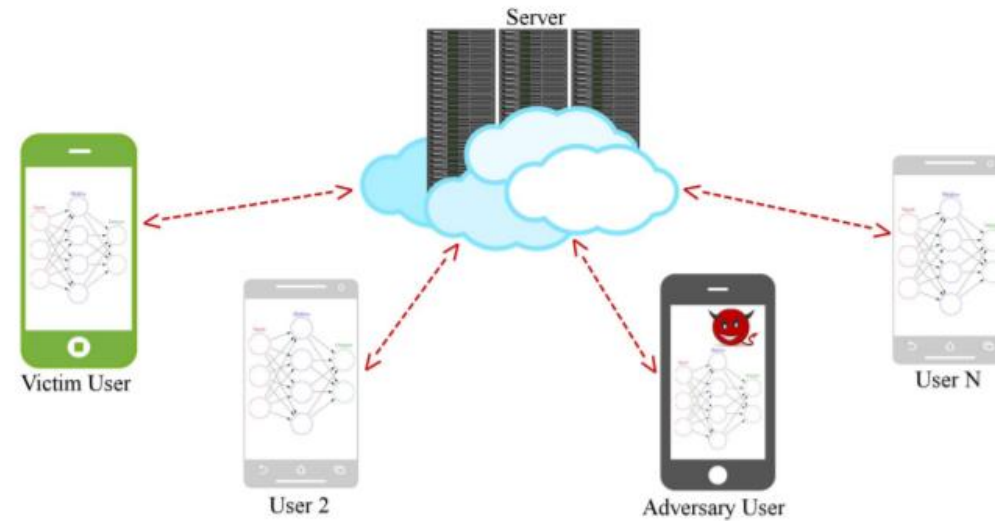**Models** [BVH+18]

**Data** [HAP17]



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.



(b) Collaborative Learning

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

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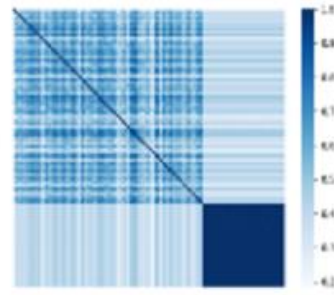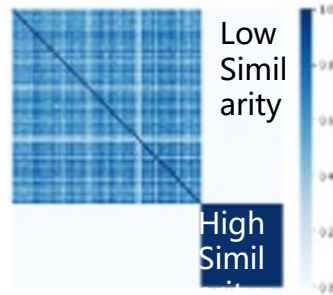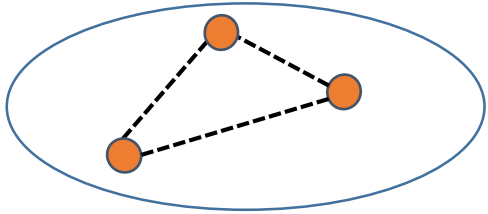
**We**Bank

# New in Federated Learning: Reliable Detection of Byzantine Colluders

Target of attack: reduce model performance

Server

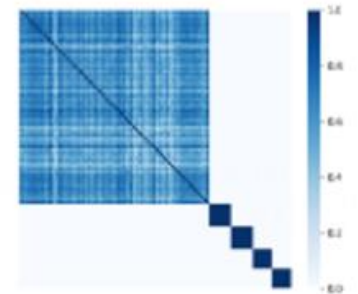Update 1  ...  Update 1

Byzantine update  ...  Byzantine update
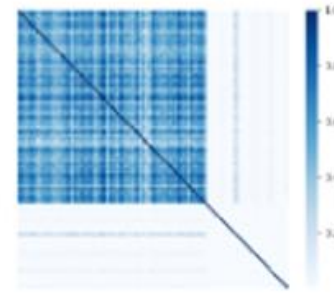
Client 1

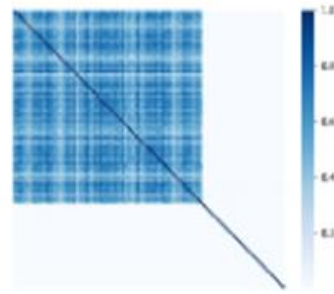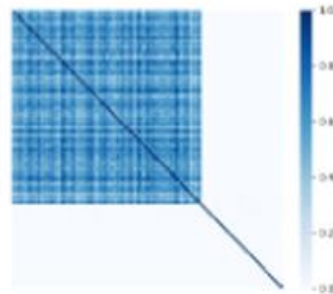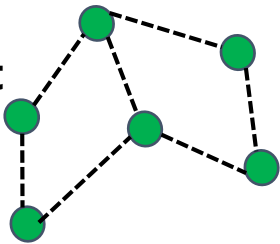Client m

Client m+1

Client K

FedSpectral: A Spectral Analysis Framework for Reliable Detection of Byzantine Colluders (submitted)

# Reliable Detection of Byzantine Colluders

**Conspiracy set 1**

**Normal client**

**Conspiracy set 2**



FedSpectral: A Spectral Analysis Framework for Reliable Detection of Byzantine Colluders (submitted)

**WeBank**

# Reliable Detection of Byzantine Colluders

- 对T轮的多个client模型更新向量，构建一个G1至GT层的**spatial-temporal图**

- 把检测看成是一个经典的图**Ncut分割**问题, 使得分割后的两个子图,
  相互间的连接最弱（最不相似），子图内部节点间的连接最强（最相似）

- 利用谱分析的eigengap heuristic来估计**共谋组的个数**，从而提升共谋检测的精确度**, 及联邦模型的性能**



**Ncut Goal:**

$$\min_{(B_1 \cup \cdots \cup B_C) = V} \sum_{t=1}^{T} \sum_{i=1}^{k} \frac{W^t(B_i, \overline{B_i})}{Vol^t(B_i)},$$

FedSpectral: A Spectral Analysis Framework for Reliable Detection of Byzantine Colluders (submitted）

# Reliable Detection of Byzantine Colluders

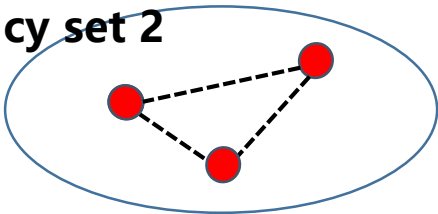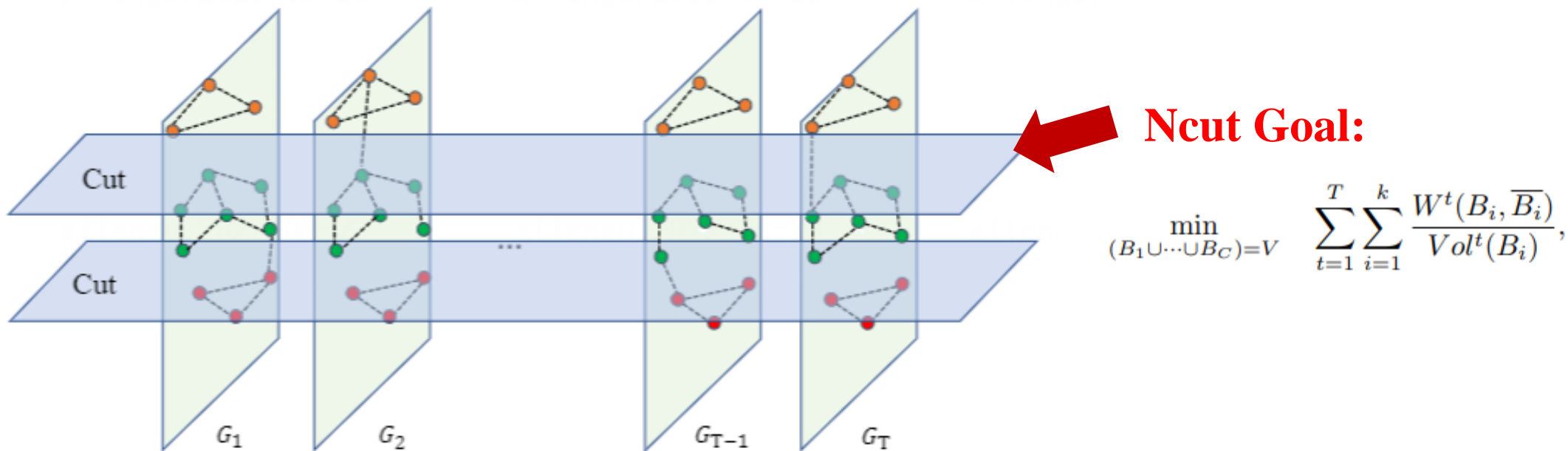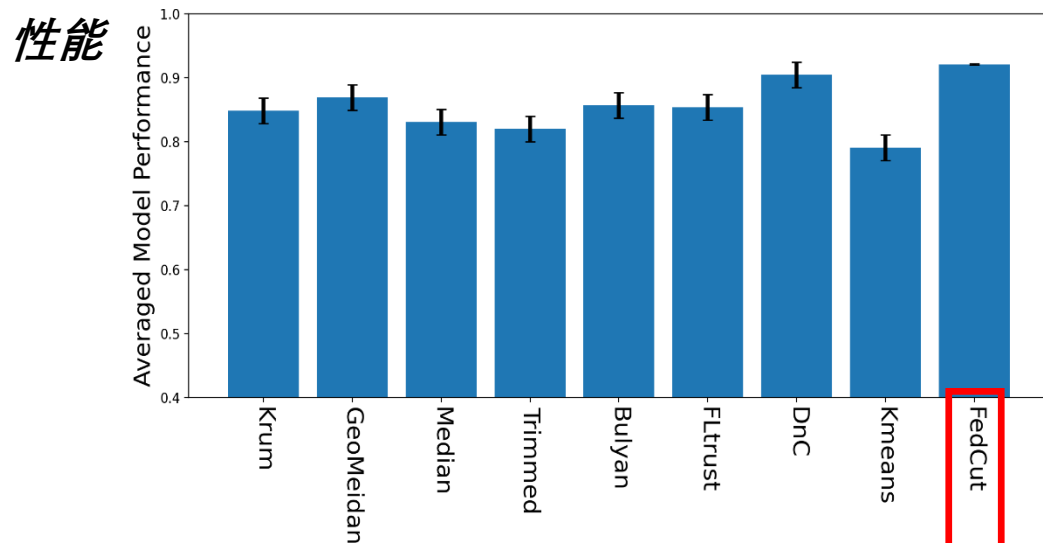*Federated Learning's Mean Performance*:

Best overall （>92%）

Worst Performance：

Federated Learning still best （>87%）
Variance small （≈5%）

性能

性能



注：图中error bar代表**不同场景多次实验**平均性能 的**方差**
（FedCut比它8种防御方法的方差更小，更鲁棒）

注：图中error bar代表**不同场景多次实验**最差性能 的**方差**
（FedCut比它8种防御方法的方差更小，更鲁棒）

FedSpectral: A Spectral Analysis Framework for Reliable Detection of Byzantine Colluders (submitted）

# 'Federated Learning' Standards

## Published

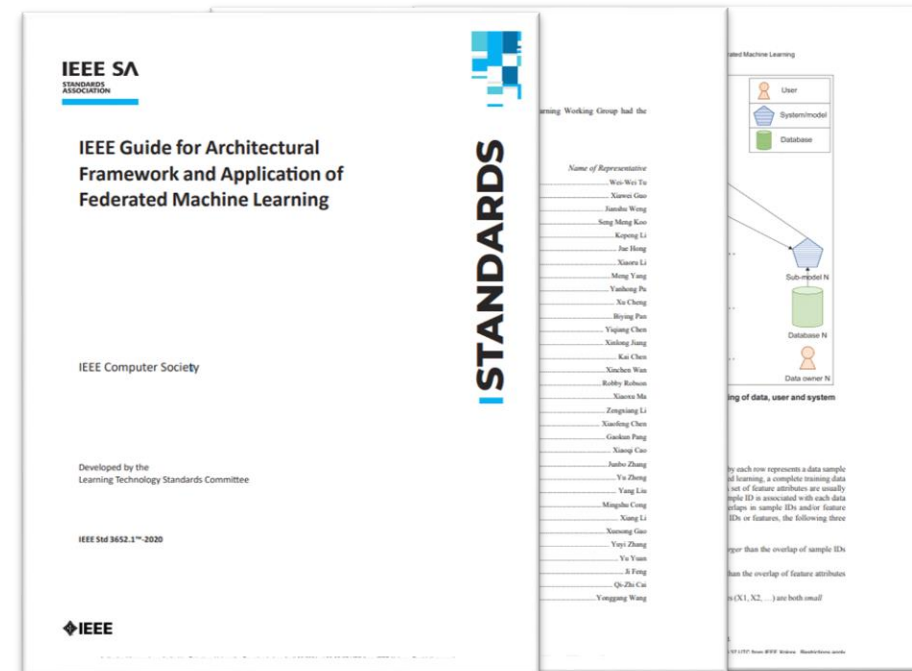IEEE P3652.1 《**Guide for Architectural Framework and Application of Federated Machine Learning**》

- **March, 2021**，**first federated learning standards in the world**，**20 participants, 6 meetings, 10 scenarios**

China Communications Standards 《**Technical Requirements and Test Methods for Federated-Learning-based Data Products** 》

- **July 2020**，**17** participants
- Scheduling management, testing etc.

## Ongoing

- Financial Industry 《 **The Financial Applications of Federated Learning and its Connectivity Specifications** 》

- Telecom Industry (CCSA-TC1/TF1)：《**The Technical Requirements and Test Methods of Safety Assessment for Federated Learning**》《**The Interoperability Technical Requirements of the Cross-Framework in Federated Learning**》

- Intra-enterprise (CCSA-T601)：《**The Interoperability Standards of the Cross-Platform in Federated Learning** 》



**WeBank** 微众银行    ◆**IEEE**    **Baidu** 百度

创新工场 SINOVATION VENTURES    **Hisense**    **4Paradigm** 第四范式

**CLUSTAR** 星云    腾讯云    京东城市

# Federated Learning Open Source Platform: FATE

**FATE**
- I. Industrial Scale
- II. Government certified
- III. Efficient, Safe, Easy to use

**Design**
- I. Many algorithms in machine learning
- II. Many privacy-preserving computation operators including MPC

**Participants include**



**Domestic Patent 427**

**International PCT  26**

**1000+** enterprises，**400+** Universities

**4600** GitHub Stars

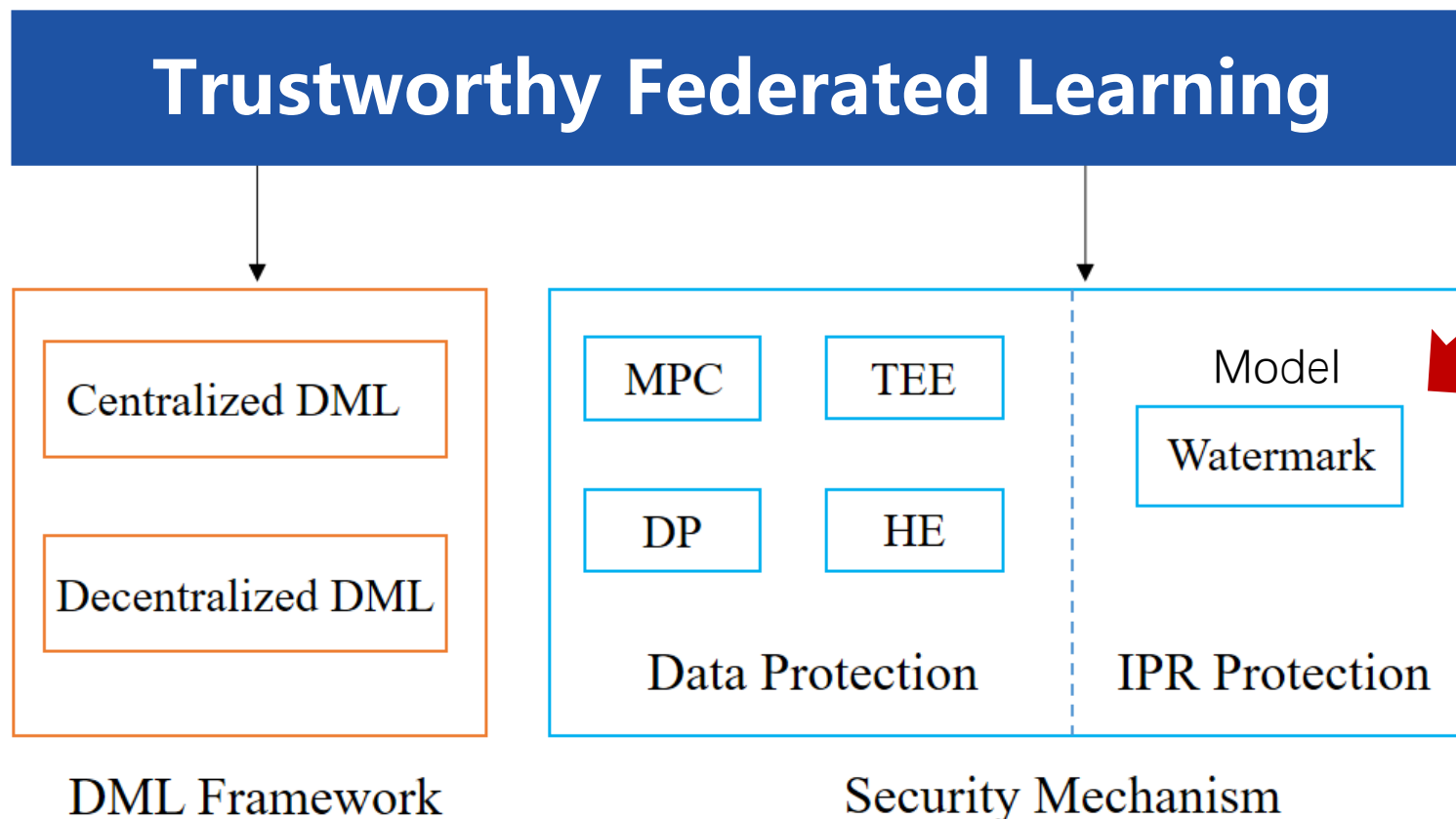19 organizations joined the FATE federated learning eco-system, including Tencent cloud, Union Pay, etc.

https://github.com/FederatedAI/

# Copyright Protection for Models (FedIPR)

- Privacy Protection

- Model Protection

- **Model IP Right**

- Explainable

- Open Source

**Trustworthy Federated Learning**

Centralized DML

Decentralized DML

**DML Framework**

| MPC | TEE |
|-----|-----|
| DP | HE |

Data Protection

Model

Watermark

IPR Protection

**Security Mechanism**

# Trustworthy Federated Learning: Protection of Model IP Rights
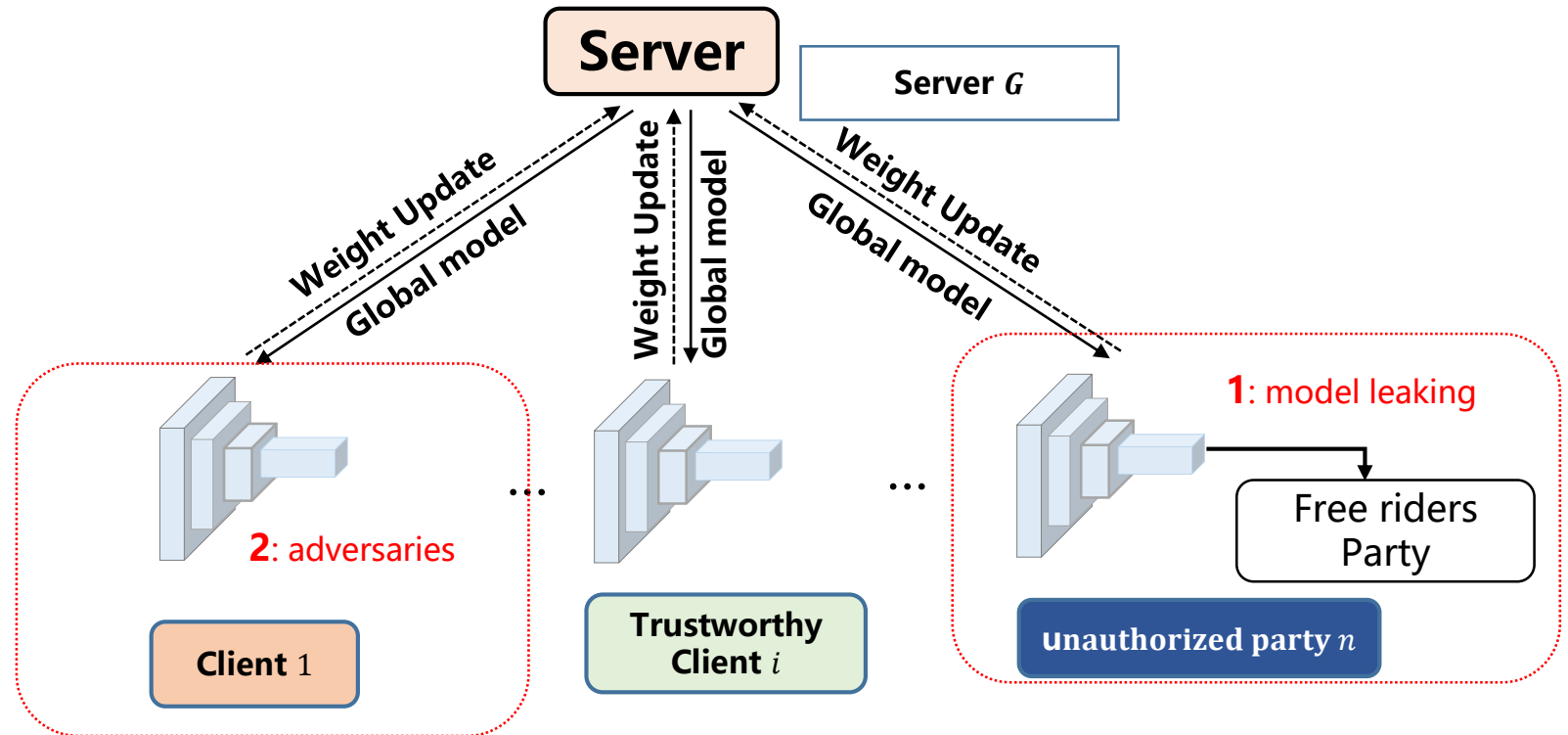
**model stealing attack [17]:**
Tramer et al.  proposed that adversaries can steal the (deployed) victim model with technical methods.

**freerider attack [18]:**
- Fraboni et al.  demonstrated the possibility that freeriders join in federated learning only for plagiarizing the valuable trained models.
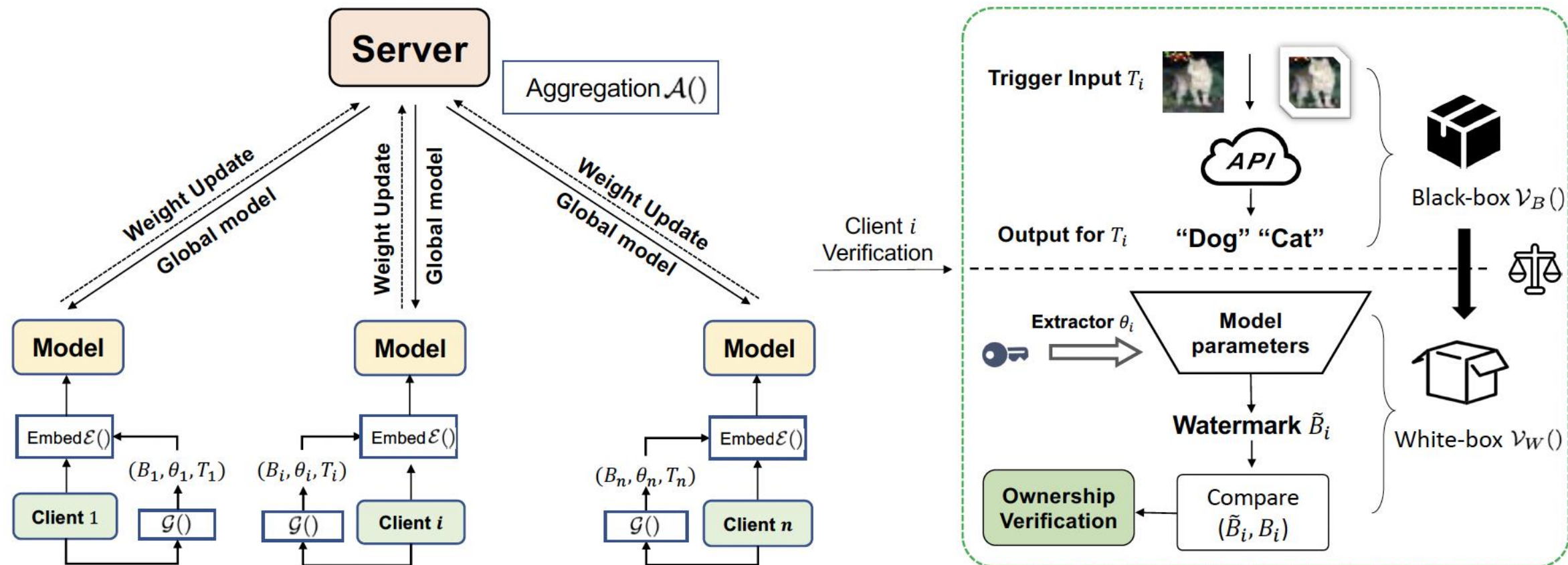
**non-technical infringement:**
- unauthorized parties may plagiarize the federated model with non-technical methods.

# Trustworthy Federated Learning: Protection of Model IP Rights

## Private watermarks

# Trustworthy Federated Learning: Protection of Model IP Rights
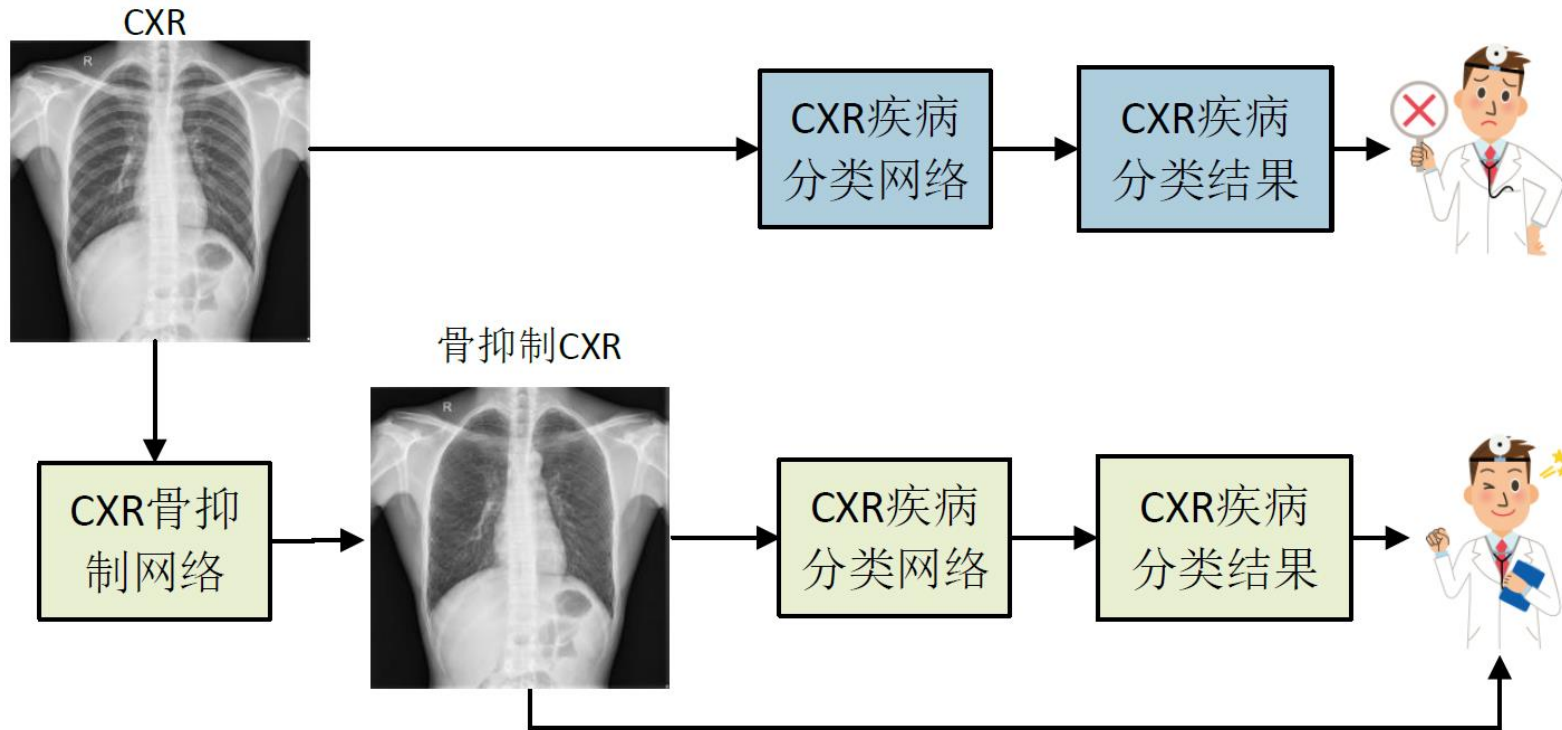


**Model tracing?**

Watermarks on data —Train→ Model —→ trustworthy

Watermarking

**Assessment functions**

# Explainability of AI Models

- A New Book



CXR

CXR骨抑制网络

骨抑制CXR

| CXR疾病分类网络 | → | CXR疾病分类结果 | ❌ |

| CXR疾病分类网络 | → | CXR疾病分类结果 |

电子工业出版社 2022年5月出版



12位人工智能领域顶级名家的扛鼎之作

阐述可解释AI研究的问题和方法
详尽展示其广泛应用和积极作用

可解释人工智能导论
Introduction to Explainable Artificial Intelligence

京东包邮
下单立减63元

全彩印刷

本书作者

杨 强　范力欣　朱 军　陈一昕　张拳石　朱松纯

陶大程　崔 鹏　周少华　刘 琦　黄萱菁　张永锋

# Figure 1: Hype Cycle for Privacy, 2021



Source: Gartner (July 2021)

Qiang Yang 202212

38

Gartner.

# Thank You