

# 联邦学习技术应用 创新探索

周旭华 博士

中国电信研究院 安全技术研究所

2022年07月09日



# 目录 CONTENT



## 01 联邦学习简介

- 背景介绍
- 联邦学习概念
- 联邦学习分类
- 与其他隐私保护计算技术

## 02 联邦学习技术创新探索

- 创新探索一：不同技术架构灵活适应不同的业务场景需求
- 创新探索二：抗数据污染或恶意窃取的新数据检测方法
- 创新探索三：抗成员推断攻击的联邦线性模型在线推理
- 创新探索四：针对纵向联邦学习的异步优化方法
- 创新探索五：参与方在联邦学习系统中对模型的价值贡献



中国电信  
CHINA TELECOM



# 01 联邦学习简介

- 背景介绍
- 联邦学习概念
- 联邦学习分类
- 与其他隐私保护计算技术



## DATA

## “数据孤岛”现象普遍存在



彼此独立的 数据孤岛

政府数据孤岛群

运营商数据孤岛群

行业数据孤岛群

# 联邦学习概念

Federated Learning Definition



中国电信  
CHINA TELECOM

DataFun.

**联邦学习**或联邦机器学习，是实现在本地原始数据不出平台的情况下，通过对中间加密数据的流通和处理来完成多方联合的机器学习训练和预测。其设计目标是在保障大数据联合价值开发时的数据安全、保护终端数据和个人数据隐私、保证合法合规的前提下，在参与多方或多计算结点之间开展高效率的机器学习，实现数据**可用不可见**。

## ◆ 数据隐私保护：

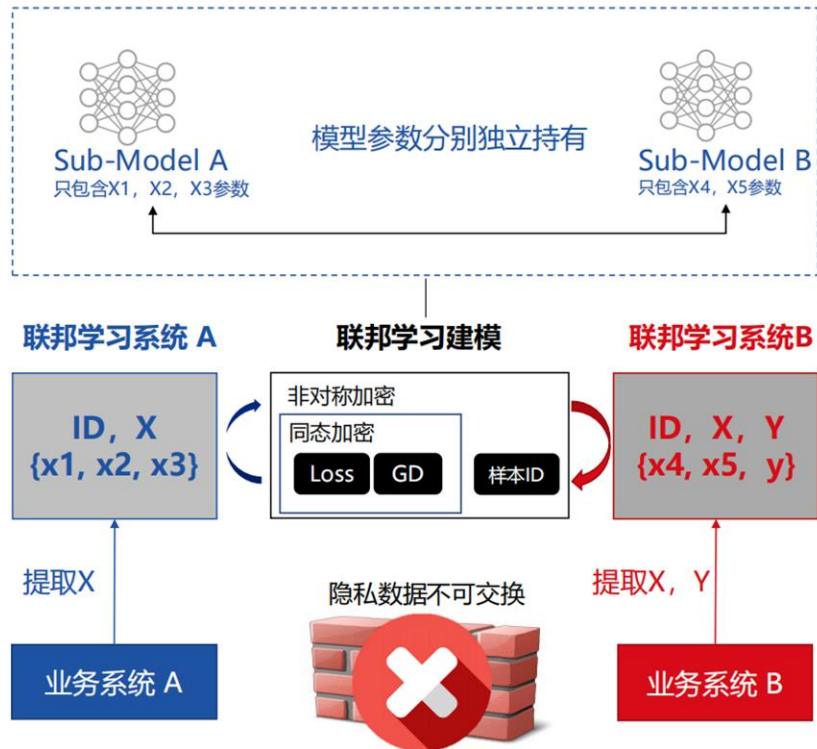
- ✓ 建模样本ID差集不向对方泄露
- ✓ 任何底层X, Y数据不向对方泄露

## ◆ 模型参数保护：

- ✓ 分别持有，联合使用

## ◆ 结果：

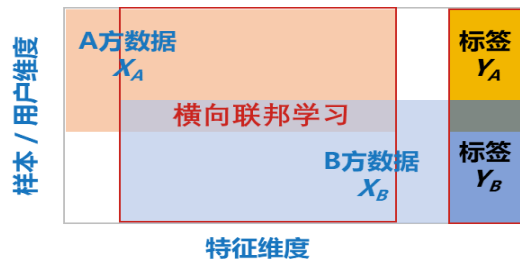
- ✓ A方有A模型
- ✓ B方有B模型
- ✓ A和B模型都比单独建模好



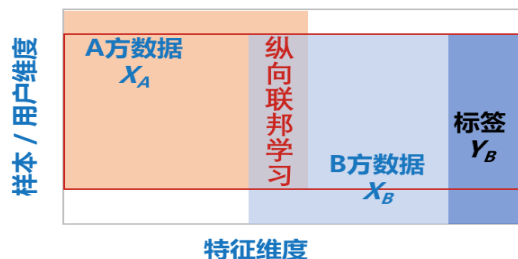
# 联邦学习分类

Federated Learning Classification

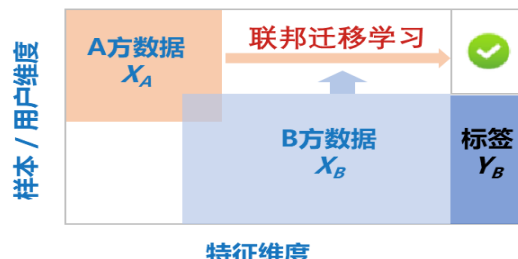
- **横向联邦学习** ➡ 可简单理解为建模算法在**不同分段数据**上进行训练和合并。
- **纵向联邦学习** ➡ 可简单理解为将**建模算法拆分**为不同模块，分别在各数据提供方进行训练，期间需要中间参数交互。



横向联邦学习



纵向联邦学习



联邦迁移学习

适用场景	同构（相同或相近领域企业合作）	异构（不同领域企业合作）	异构（不同领域企业合作）
数据特征	特征重叠较多、样本重叠较少	特征重叠较少、样本重叠较多	特征与样本重叠皆较少
工程成熟度	较成熟，大型企业广泛应用	已进入发展应用阶段，部分企业已开始应用	应用较少

# 与其他隐私保护计算技术

With Other Privacy-preserving Computing Technologies

## 联邦学习特点：



针对传统集中模型训练存在泄露数据隐私问题而提出，增加交互安全设计



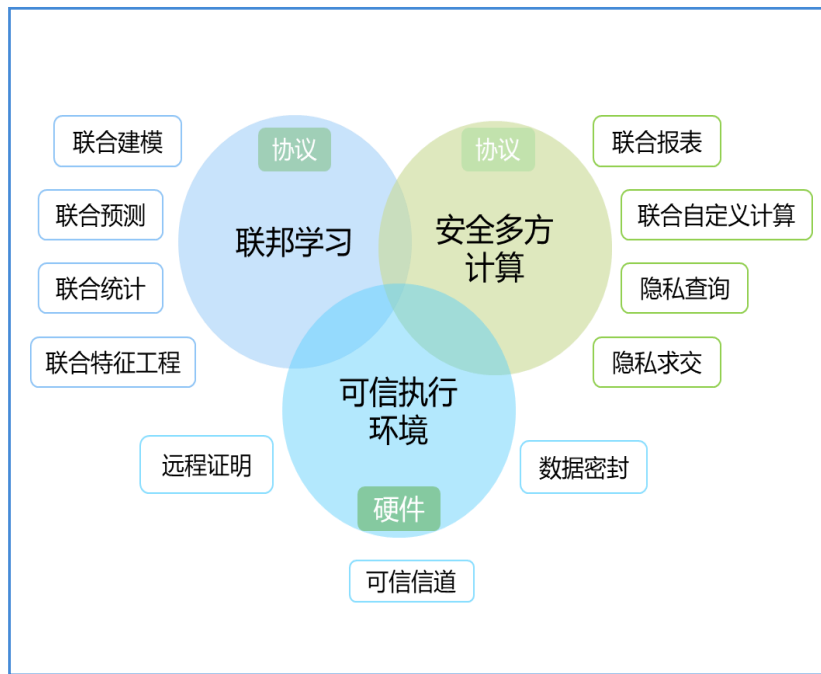
强调“数据不出平台”，核心理念是“数据不动模型动”



模型性能接近或几乎无损



提供技术框架，容易与其他隐私保护计算技术结合



多种隐私保护计算技术交叉应用



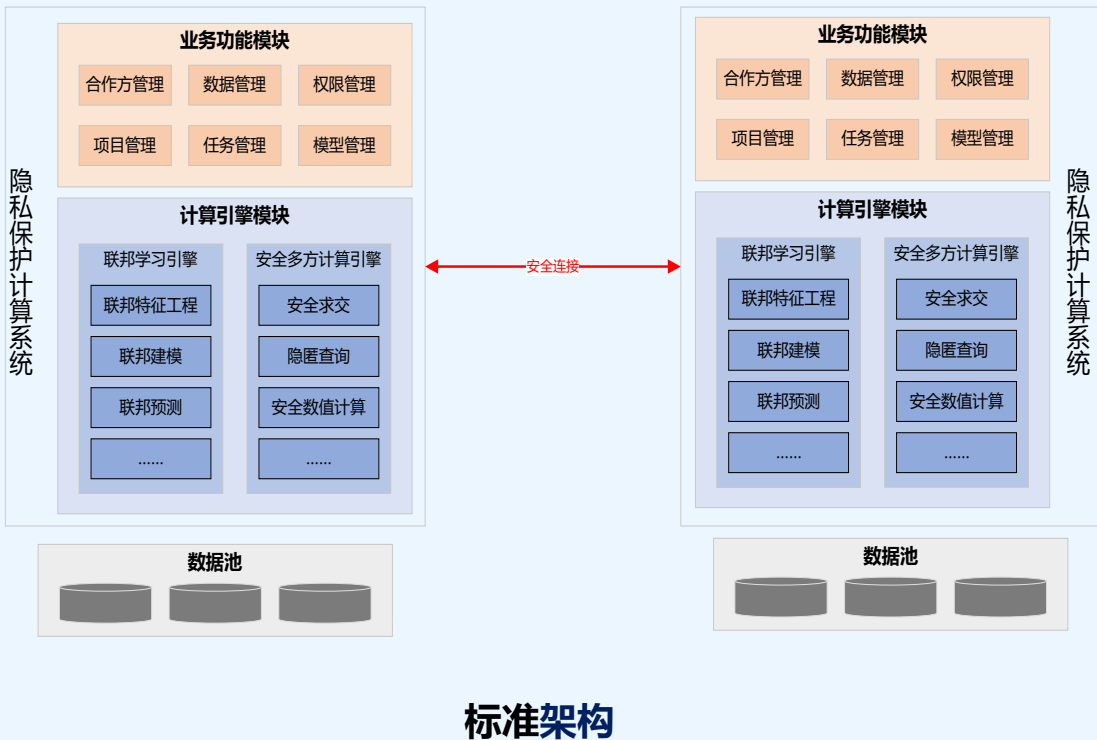
## 02 技术创新探索

- 创新探索一：不同系统架构灵活适应不同的业务场景需求
- 创新探索二：抗数据污染或恶意窃取的新数据检测方法
- 创新探索三：抗成员推断攻击的联邦线性模型在线推理
- 创新探索四：针对纵向联邦学习的异步优化方法
- 创新探索五：参与方在联邦学习系统中对模型的价值贡献





## 创新探索一：标准架构、交易中心架构灵活适应不同的业务场景需求



### 标准版特点一

可满足中国电信作为数据提供方向多行业提供电信数据开发数据价值的需要



### 标准版特点二

解决了现网应用遇到的亿级大数据量、兆级低带宽高延时、网络不直达问题

## 创新探索一：标准架构、交易中心架构灵活适应不同的业务场景需求



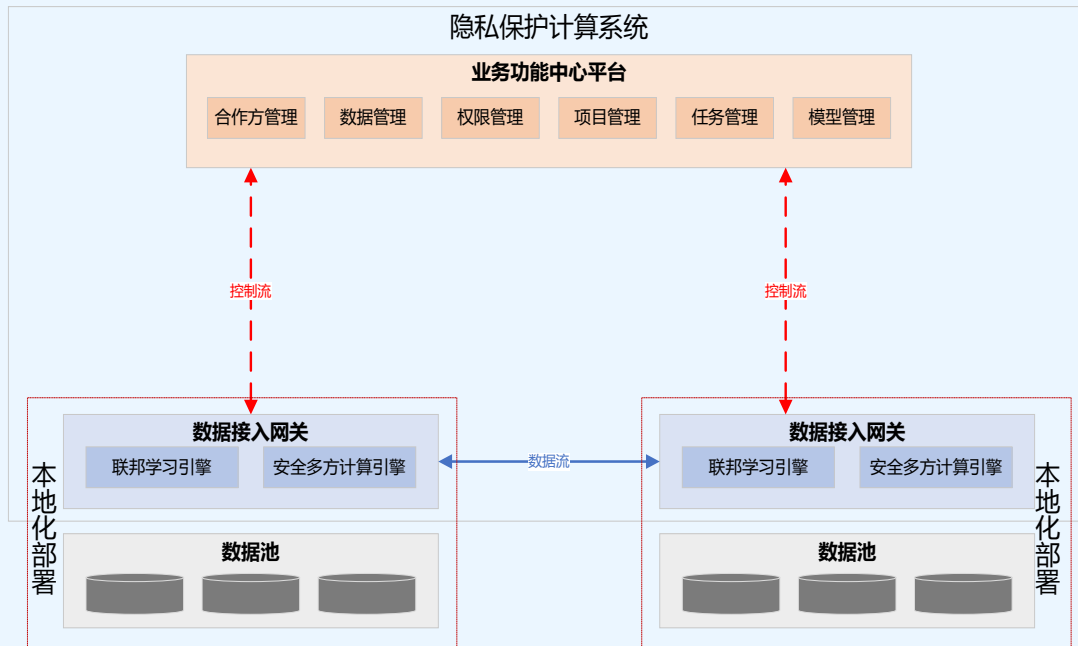
### 交易中心版特点一

实现管理模块与计算模块相分离



### 交易中心版特点二

保证数据不出各方管理域的前提下，做到统一管理、统一入口



交易中心架构

## 创新探索二：纵向联邦学习场景的数据污染的新数据检测方法

中国电信  
CHINA TELECOM

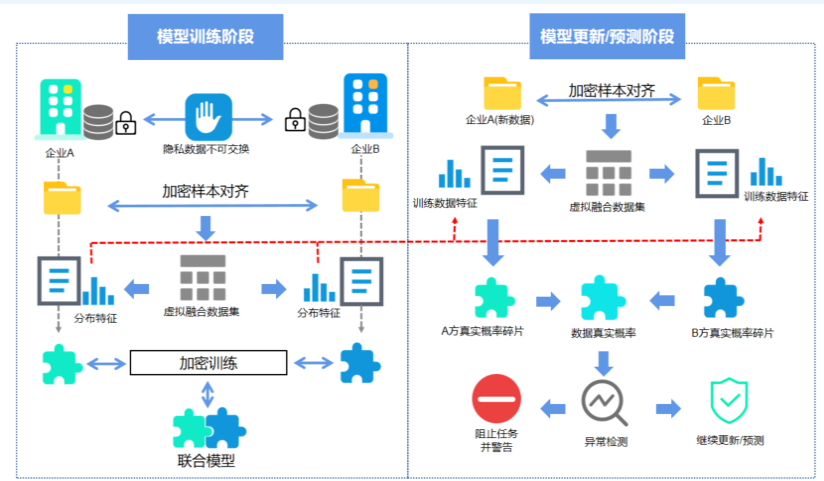
DataFun.

存储基于一批有效训练数据得到的各特征的概率分布特征

根据各特征的概率分布特征设置正常出现概率阈值 $\epsilon$

计算单条数据 $x$ 的出现概率 $P(x)$

将出现概率 $P(x)$ 与概率阈值 $\epsilon$ 进行比较来判定是否异常



联邦学习数据污染检验技术探索



### 存在问题

在联邦学习框架下参与方不能直接访问他方的数据，污染数据或恶意数据更有可能发生，使模型失效。



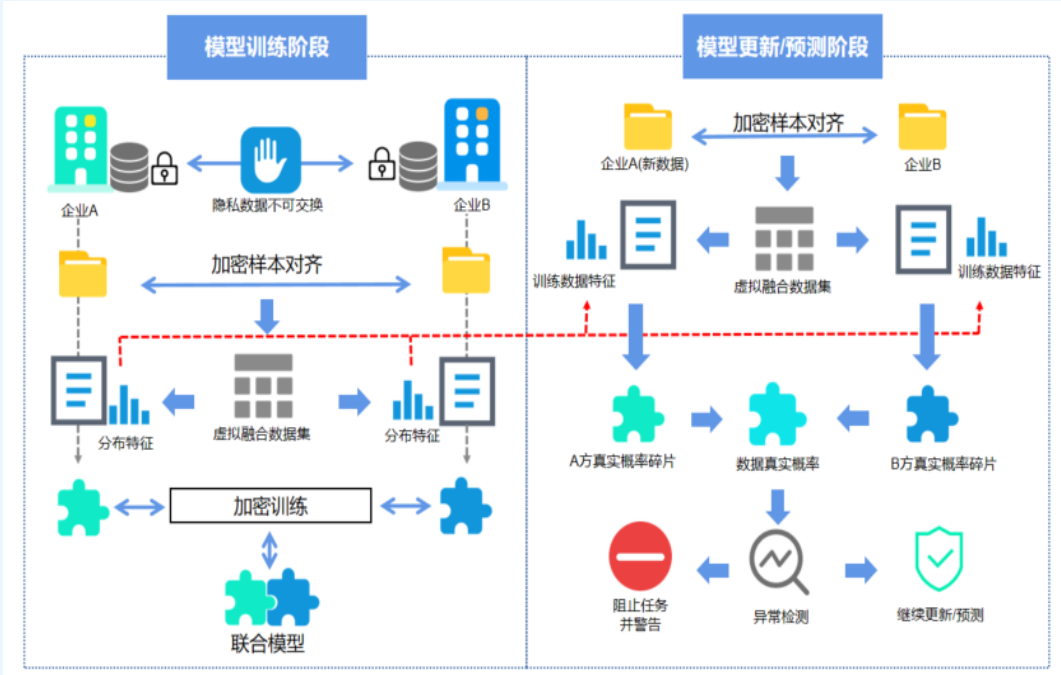
### 数据污染检方案

通过训练数据特性的提取与存储以及针对其不同特征的分类分析来计算参与方新提供数据是有效数据的概率。

## 创新探索二：纵向联邦学习场景的数据污染的新数据检测方法

中国电信  
CHINA TELECOM

DataFun.



联邦学习数据污染检验技术探索



### 存在问题

在联邦学习框架下参与方不能直接访问他方的数据，污染数据或恶意数据更有可能发生，使模型失效。



### 数据污染检方案

通过训练数据特性的提取与存储以及针对其不同特征的分类分析来计算参与方新提供数据是有效数据的概率。

## 创新探索三：抗成员推断攻击的联邦线性模型在线推理



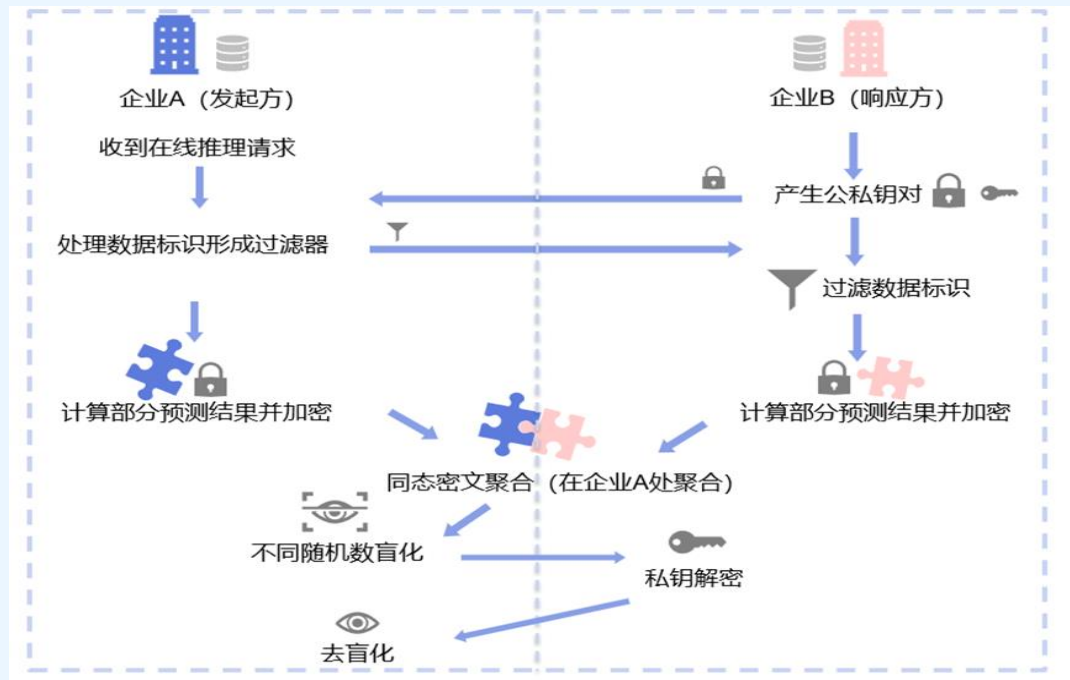
### 存在问题

联邦模型通过联邦在线推理提供数据预测功能，传统过程中使用明文方式进行，存在数据泄露和用户隐私信息问题。



### 安全在线推理方案

利用过滤器、同态加密算法和随机数乘法盲化法保护发起方的请求无法被响应方精确获悉，从而抵抗成员推断攻击。



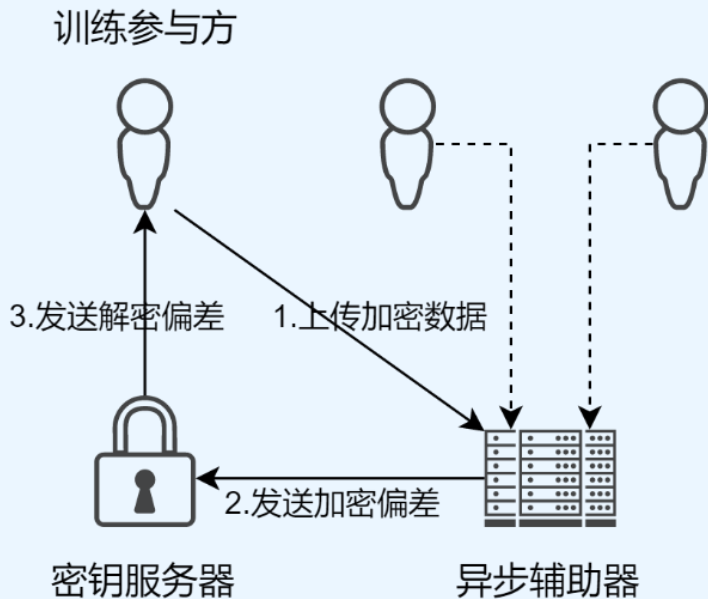
### 联邦在线推理技术探索

## 创新探索四：针对纵向联邦学习的异步优化方法



中国电信  
CHINA TELECOM

DataFun.



联邦异步加速技术探索



### 存在问题

每个参与方数据量、计算速度、网络延迟都不一致，训练时间由处理最慢的参与方决定，从而形成木桶效应，影响训练效率。



### 联邦异步加速方案

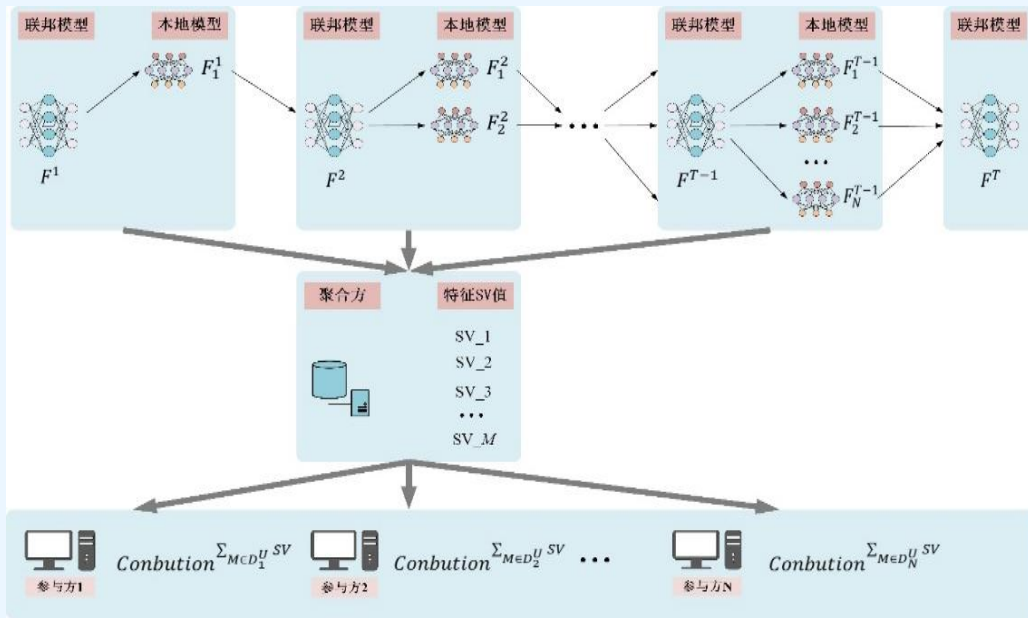
引入存储参与方特征计算值的缓存单元，打破参与方编码模型训练时的相互依赖，实现各参与方编码模型之间的异步训练。

## 创新探索五：参与方在联邦学习系统中对模型的价值贡献



中国电信  
CHINA TELECOM

DataFun.



联邦模型贡献量评估



### 存在问题

当前联邦学习系统中忽视各参与方数据对联邦模型增益的贡献差异，无法推动跨域跨行业数据共享的良性循环。



### 贡献量评估方案

探索提出基于SV理论的贡献量评估方法，通过考虑特征重要性来反映各方数据在联邦学习系统对模型的价值贡献。



# 非常感谢您的观看

---

