# A Quasi-Newton Method Based Vertical Federated Learning Framework for Logistic Regression

Kai Yang♠♣, Tao Fan♠, Tianjian Chen♠, Yuanming Shi♣, and Qiang Yang♦

♠WeBank
♣ShanghaiTech University
♦Hong Kong University of Science and Technology

WeBank 微众银行

上海科技大学 ShanghaiTech University
香港科技大學 THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY

## Privacy-Preserving Collaborative Machine Learning

**Main concerns:** data privacy and security for building a machine learning model

Training machine learning model at cloud data center: risks of data breach and violation of data protection laws and regulations (e.g., GDPR by the European Union)

**Federated learning:** an emerging frontier field on privacy-preserving collaborative machine learning while leaving data instances at their providers locally [1]

- **horizontal federated learning structure:** each node has a subset of data instances with complete data attributes
- **vertical federated learning structure:** each node holds a disjoint subset of attributes for all data instances

**Communication challenge:** one of the main bottlenecks in federated learning due to the much worse network conditions than the cloud center

## Vertical Federated Learning

Vertical federated learning framework: joint computation and communication design for different ML models (e.g., logistic regression, boosting-tree, etc.)

**Design target:**
- preserving data privacy
- low communication costs

**State-of-the-art:** SGD proposed in [2] based on Taylor expansion and additive homomorphic encryption

**Challenges:**
- high communication costs due to low convergence rate of SGD
- high computational costs of second-order methods

**Proposal:** computationally efficient **quasi-Newton method** to improve the convergence rate without introducing much additional communication costs at each iteration

**Quasi-Newton methods:**
- **L-BFGS:** high communication costs for transmitting inverse Hessian matrix
- **stochastic quasi-newton method in [Schraudolph, et al., 2007]:** unstable estimation for the inverse Hessian matrix with small batch sizes
- **stochastic L-BFGS in [Moritz, et al., 2016]:** requiring computing the full gradient for approximating inverse Hessian matrix, which doubles the average communication cost at each iteration
- **stochastic quasi-newton method in [3]:** updating the approximated inverse Hessian matrix every $L$ iterations based on a gradient-like vector

We develop a communication efficient vertical federated learning framework based on the stochastic quasi-Newton method proposed in [3].

## Problem Statement: Vertical Logistic Regression

**Problem setting** of logistic regression:
- $X \in \mathbb{R}^{n \times T}$: data set consisting of $T$ data samples and each instance has $n$ features
- $y \in \{-1, +1\}^T$: labels of $X$
- $w$: model parameters
- $x_i$: $i$-th data instance
- $y_i$: label of $x_i$
- $l(w; x_i, y_i) = \log(1 + \exp(y_i w^\mathsf{T} x_i))$: negative log-likelihood loss
- **Target:**
$$\underset{w \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{T}\sum_{i}^{T} l(w; x_i, y_i), \qquad (1)$$

**System setting** of vertically federated learning for logistic regression: each party holds a disjoint subset of data features over a common sample IDs

- A and B: two honest-but-curious private parties
- A: the **host** data provider with only features ($X^A \in \mathbb{R}^{n_A \times T}$)
- B: the **guest** data provider with features ($X^B \in \mathbb{R}^{n_B \times T}$) and labels $y \in \{-1, +1\}^T$

Vertically partitioned model parameters: party $A$ and party $B$ hold the model parameters corresponding to their features respectively, i.e., $w = (w^A \in \mathbb{R}^{n_A}, w^B \in \mathbb{R}^{n_B})$

**Additively homomorphic encryption** for exchanging encrypted intermediate values: e.g., Paillier.

- Encryption: $[\![u]\!] + [\![v]\!] = [\![u+v]\!]$, $v \cdot [\![u]\!] = [\![vu]\!]$ where $[\![\cdot]\!]$ is the encryption operation
- Decryption: requiring a third party called the **coordinator**

**Taylor loss:** $l(w; x_i, y_i) \approx \log 2 - \frac{1}{2} y_i w^\mathsf{T} x_i + \frac{1}{8}(w^\mathsf{T} x_i)^2$ second-order Taylor approximation for loss function.

## Our Work: A Quasi-Newton Method for Vertical Logistic Regression

### Quasi-Newton Method Based Vertical Federated Learning

**Target:** reducing communication rounds without increasing much communication bandwidth at per round with quasi-Newton method

**Key ideas:**
- **curvature information** $H$: estimated inverse Hessian matrix
- Update of quasi-Newton method: $w \leftarrow w - \eta H g$
- **Subsampled method** for curvature estimation [3]: updating $H$ every $L$ iterations to reduce the communication overhead as well as improve the stability of quasi-Newton algorithm

**Gradient of Taylor loss:** $\nabla l(w; x_i, y_i) \approx (\frac{1}{4} w^\mathsf{T} x_i - \frac{1}{2} y_i) x_i$    **Hessian of Taylor loss:** $\nabla^2 l(w; x_i, y_i) \approx \frac{1}{4} x_i x_i^\mathsf{T}$

**Proposed framework:**

- **Computing Loss and Gradient at Party A&B:**

At each iteration, choose a mini-batch of data instances: $\mathcal{S} \subseteq \{1, \cdots, T\}$ is the index set.

- loss and gradient loss, $g$: loss $= F(w) = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} l(w; x_i, y_i)$, $g = \nabla F(w) = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \nabla l(w; x_i, y_i)$
- intermediate values $u_A, u_A^2, u_B, u_B^2, d$: $u_A = \{u_A[i] = w^{A\mathsf{T}} x_i^A : i \in \mathcal{S}\}$, $u_A^2 = \{u_A^2[i] = (w^{A\mathsf{T}} x_i^A)^2 : i \in \mathcal{S}\}$ $u_B = \{u_B[i] = w^{B\mathsf{T}} x_i^B : i \in \mathcal{S}\}$, $u_B^2 = \{u_B^2[i] = (w^{B\mathsf{T}} x_i^B)^2 : i \in \mathcal{S}\}$ $d = \{d_i = \frac{1}{4}(u_A[i] + u_B[i]) - \frac{1}{2}y_i) : i \in \mathcal{S}\}$
- **encrypted loss and gradient** $[\![\text{loss}]\!], [\![g]\!]$:
$[\![\text{loss}]\!] \approx \frac{1}{|\mathcal{S}|}\sum_{i \in \mathcal{S}}[\![\log 2]\!] - \frac{1}{2}y_i([\![u_A[i]]\!] + [\![u_B[i]]\!]) + \frac{1}{8}([\![u_A^2[i]]\!] + 2u_B[i][\![u_A[i]]\!] + [\![u_B^2[i]]\!])$, $[\![g]\!] \approx \frac{1}{|\mathcal{S}|}\sum_{i \in \mathcal{S}}[\![d_i]\!]x_i = ([\![g^A]\!], [\![g^B]\!]) = (\sum_{i \in \mathcal{S}}[\![d_i]\!]x_i^A, \sum_{i \in \mathcal{S}}[\![d_i]\!]x_i^B)$, $[\![d_i]\!] = \frac{1}{4}([\![u_A[i]]\!] + [\![u_B[i]]\!] + [\![-\frac{1}{2}y_i]\!])$.

At each iteration, by transmitting $[\![u_A]\!]$ from party A to party B, and transmitting $[\![d]\!]$ from B to A, $[\![g^A]\!]$ can be computed at party A, $[\![\text{loss}]\!]$ and $[\![g^B]\!]$ can be computed at party B privately.

- **Computing Updates for Estimating Curvature Information at Party A&B:**

Every $L$ iterations, choose a subset of data instances for estimating curvature information: $\mathcal{S}_H$. The coordinator collects encrypted $v = (v^A, v^B) \in \mathbb{R}^n$ from party A and B for updating the curvature information $H$.

- difference of average model parameters $s_t$: $s_t = \bar{w}_t - \bar{w}_{t-1} = (s_t^A, s_t^B)$, $\bar{w}_t = \sum_{i=k-L+1}^{k} w_i/L$, $\bar{w}_{t-1} = \sum_{i=k-2L+1}^{k-L} w_i/L$
- $v, h$: $v_t = \nabla^2 \hat{F}(\bar{w}_t)s_t$, where $\nabla^2 \hat{F}(\bar{w}_t) = \frac{1}{|\mathcal{S}_H|}\sum_{i \in \mathcal{S}_H} \nabla^2 l(\bar{w}_t; x_i, y_i) = \frac{1}{|\mathcal{S}_H|}\sum_{i \in \mathcal{S}_H} x_i x_i^\mathsf{T}$. $h = \{h_i = \Delta \bar{u}_A^A + \Delta \bar{u}_B^B = s_t^{A\mathsf{T}}x_i^A + s_t^{B\mathsf{T}}x_i^B, i \in \mathcal{S}_H\}$.
- Computing $[\![v]\!]$: $[\![v_t]\!] = \frac{1}{|\mathcal{S}_H|}\sum_{i \in \mathcal{S}_H}[\![h_i]\!]x_i = ([\![v_t^A]\!], [\![v_t^B]\!]) = (\frac{1}{|\mathcal{S}_H|}\sum_{i \in \mathcal{S}_H}[\![h_i]\!]x_i^A, \frac{1}{|\mathcal{S}_H|}\sum_{i \in \mathcal{S}_H}[\![h_i]\!]x_i^B)$,

Every $L$ iterations, by transmitting $[\![\Delta \bar{u}_A]\!] = \{[\![\Delta \bar{u}_i^A]\!] : i \in \mathcal{S}_H\}$ from party A to party B, and transmitting $[\![h]\!] = \{[\![h_i]\!] : i \in \mathcal{S}_H\}$ from B to A, $[\![v_t^A]\!]$ can be computed at party A and $[\![v_t^B]\!]$ can be computed at party B privately.

- **Computing Descent Direction at the Coordinator:** By decryption the coordinator obtains loss, $g$, $v$ from party A and B.

At each iteration, the coordinator should determine a descent direction $\bar{g}$ for updating $w^A$ and $w^B$: $w \leftarrow w - \bar{g} = w - \eta H g = (w^A - \bar{g}^A, w^B - \bar{g}^B)$.

Every $L$ iterations, the coordinator should also update $H$ based on the collected encrypted loss $[\![\text{loss}]\!]$, gradient $[\![g]\!]$, and $[\![v]\!]$ from party A&B.

- Initial point: $H = (v_t^\mathsf{T} s_t / v_t^\mathsf{T} v_t)I$. For $\forall j = t - M + 1, \cdots, t$, iteratively compute $H \leftarrow (I - \rho_j s_j v_j^\mathsf{T})H(I - \rho_j v_j s_j^\mathsf{T}) + \rho_j s_j s_j^\mathsf{T}$, $\rho_j = 1/(v_j^\mathsf{T} s_j)$,

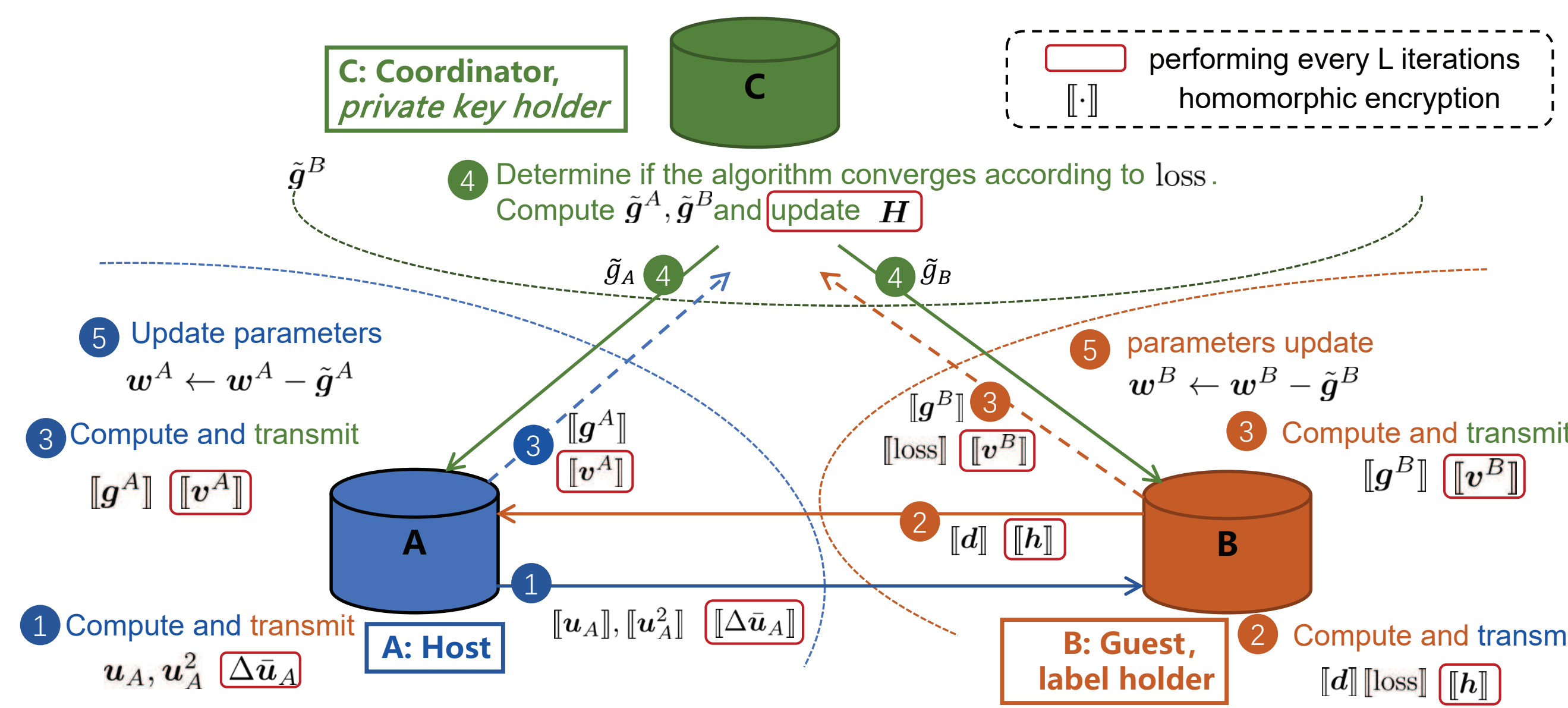The source code will be released in an upcoming version of the FATE framework [4].



Figure 1: Proposed Framework for Vertical Federated Learning

**Algorithm 1:** Proposed Framework for Vertical Federated Learning

**Input** : $w_0^A, w_0^B, M, L$
**Output**: $w^A, w^B$

1  Set $t = 0, H = I$
2  **for** each round $k = 1, \cdots,$ **do**
3    Choose a minibatch $\mathcal{S}$
4    **if** $mod(k, L) \neq 0$ **then**
5      **Party A&B:** compute $[\![\text{loss}]\!], [\![g]\!]$ as equation (3) (4)
6      **Coordinator:** $w_{k+1} = w_k - \bar{g}_k$ where $\bar{g}_k = \eta H g$
7    **else**
8      $t \leftarrow t + 1$
9      **Party A&B:** Choose a minibatch $\mathcal{S}_H$
10     compute $[\![\text{loss}]\!], [\![g]\!], [\![v_t]\!]$ as equation (3) (4) (6)
11     **Coordinator:** $w_{k+1} = w_k - \bar{g}_k$ where $\bar{g}_k = \eta H g$
12     $s_t = \sum_{i=k-L+1}^{k} \bar{g}_i/L - \sum_{i=k-2L+1}^{k-L} \bar{g}_i/L$
13     **if** $t > 1$ **then**
14       $H \leftarrow (s_t^\mathsf{T} v_t)/(v_t^\mathsf{T} v_t)I, \tilde{m} = \min\{M, t\}$
15       **for** $j = t - \tilde{m} + 1, \cdots, t$ **do**
16         $\rho_j = 1/(v_j^\mathsf{T} s_j)$
17         $H \leftarrow (I - \rho_j s_j v_j^\mathsf{T})H(I - \rho_j v_j s_j^\mathsf{T}) + \rho_j s_j s_j^\mathsf{T}$
18       **end**
19     **end**
20     $\bar{w}_t = 0$
21   **end**
22  **end**

## Communication Costs of Each Iteration

**Communication costs of SGD [2]:** $3|\mathcal{S}|$ encrypted numbers between party A and party B, and $2n$ encrypted numbers between party A&B and the coordinator.

**Communication costs of proposed framework:** $3|\mathcal{S}| + 2|\mathcal{S}_H|/L$ encrypted numbers between party A and party B, and $(2 + 1/L)n$ encrypted numbers between party A&B and the coordinator. By choosing $|\mathcal{S}_H| \leq |\mathcal{S}|$, the presented quasi-Newton method introduces no more than $1/L$ additional communication costs at per communication round compared with [2].

## Experiments

Numerical experiments on two credit scoring data sets by setting $\mathcal{S}_H = \mathcal{S}$ and $L = 4$:
- **Credit 1:** 30000 data instances and $n = 25$ attributes.
- **Credit 2:** 150000 data instances and 10 attributes.

Table 1: Numerical Results on Two Public Data Sets

| Batch Size | Method | Credit 1 | | | Credit 2 | | |
|---|---|---|---|---|---|---|---|
| | | Epochs | Loss | AUC | Epochs | Loss | AUC |
| 1000 | SGD | 12 | 0.496218 | 0.7224 | 12 | 0.314555 | 0.7033 |
| | Proposed | 3 | 0.496600 | 0.7222 | 4 | 0.314643 | 0.7061 |
| 3000 | SGD | 18 | 0.496194 | 0.7219 | 14 | 0.314648 | 0.6982 |
| | Proposed | 12 | 0.496317 | 0.7225 | 6 | 0.314490 | 0.7077 |

## Conclusions

- Addressing the communication challenge in vertical federated learning for logistic regression.
- A quasi-Newton framework to reduce the number of communication rounds without introducing much additional communication costs at each round.
- Computing an encrypted gradient and an additional vector every $L$ iterations for updating the curvature information with additively homomorphic encryption.
- Advantages demonstrated via numerical experiments.

## References

[1] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong, "Federated machine learning: Concept and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):12, 2019.
[2] Stephen Hardy, Wilko Henecka, Hamish Ivey-Law, Richard Nock, Giorgio Patrini, Guillaume Smith, and Brian Thorne, "Private federated learning on vertically partitioned data via entity resolution and additively homomorphic encryption," *arXiv preprint arXiv:1711.10677*, 2017.
[3] Richard H Byrd, Samantha L Hansen, Jorge Nocedal, and Yoram Singer, "A stochastic quasi-newton method for large-scale optimization," *SIAM Journal on Optimization*, 26(2):1008–1031, 2016.
[4] WeBank. FATE: An industrial grade federated learning framework. https://fate.fedai.org, 2018.

## Contact Information

- Kai Yang: yangkai@shanghaitech.edu.cn
- Tao Fan: dylanfan@webank.com

FATE