# Fortifying Federated Learning Towards Trustworthiness via Auditable Data Valuation and Verifiable Client Contribution

K Naveen Kumar[1], Ranjeet Ranjan Jha[2], C Krishna Mohan[1], and Ravindra Babu Tallamraju[3]

[1]IIT Hyderabad, India [2]IIT Patna, India [3]Sahaj AI Software Pvt Ltd., India

## Problem Statement

- Federated Learning (FL) facilitates clients to collaborate on training a shared machine learning model without **exposing individual private data.**
- Requires **auditability** and **verifiability** to ensure **local data integrity** and **trustworthy client updates**.
- Existing FL frameworks rely on a server for **client selection** and **aggregation**, creating a <span style="color:red">single point of failure</span> and <span style="color:red">privacy risks</span>.

## Proposed Solution: FAVD (Figure 1)

- **Federated auditable and verifiable data valuation (FAVD)** ensures auditability and verifiability of client contributions without a central authority.
- Utilizes **local data density functions** to construct a global density function for transparent and effective data valuation before training.
- **Gaussian noise** is added to shared density functions to mitigate privacy risks.

## Key Contributions

- **Privacy-preserving data valuation** independent of any predefined training algorithm.
- Ensures **benign updates** even in the presence of **malicious data.**
- Theoretical analysis on **convergence, auditability, verifiability**, and **resilience** against **data poisoning threats.**
- Improves fairness by **distinguishing** high-impact contributions from low-quality data.
- **Comprehensive evaluation on five benchmark datasets with various models:** Covid-chestxray (ResNet50), Camelyon17 (DenseNet121), HAM10000 (FL-FixCaps), CIFAR10 (ResNet18), CIFAR100 (ResNet50).
- **Introduces novel metrics for auditable data valuation**
  - Malicious Sample Detection Rate **(MSDR)**
  - Benign Misclassification Rate **(BMR)**
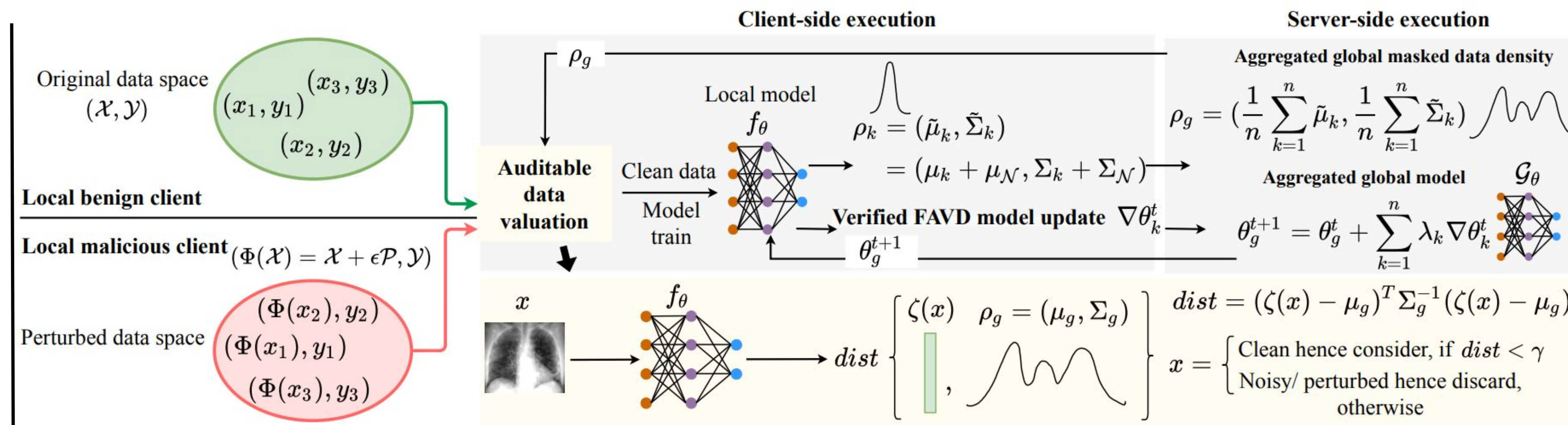  - Client Contribution Consistency **(CCC)** Score



Figure 1. Overview of the FAVD-integrated FL system, with client-side and server-side execution.

**Algorithm 1** Standard FL with **our FAVD framework**

**Input:** Global model $\mathcal{G}_{\theta,t}$, local data $\mathcal{D}_k = (\mathcal{X}_k, \mathcal{Y}_k)$, privacy noise parameters $\rho_\mathcal{N} = (\mu_\mathcal{N}, \Sigma_\mathcal{N})$, auditing parameter $\gamma$

**Output:** Global test accuracy $\mathcal{A}_\mathcal{G}$

1: **Client execution** $(\theta_g^t, \rho_g = (\mu_g, \Sigma_g))$:
2: **for** each client $k = 1$ **to** $n$ **do**
3:    Initialize the local model $\theta_k^t \leftarrow \theta_g^t$
4:    **if** client $k$ is malicious **then**
5:      $\Phi(\mathcal{X}_k) \leftarrow \mathcal{X}_k + \epsilon\mathcal{P}$
6:      $\mathcal{D}_k(\mathcal{X}_k) \leftarrow \tilde{\mathcal{D}}_k(\Phi(\mathcal{X}_k))$
7:    $\mathcal{D}_k^\mathcal{F}, \rho_k \leftarrow \text{FAVD}(\mathcal{D}_k, \gamma, \rho_\mathcal{N}, f_{\theta,k})$
8:    **for** $b = 1$ **to** batches $\in \mathcal{D}_k^\mathcal{F} = (\mathcal{X}_k^\mathcal{F}, \mathcal{Y}_k)$ **do**
9:      $\{\mathcal{Q}_{b,k}\} \leftarrow \{\sigma(f_{\theta,k}(\mathcal{X}_k^\mathcal{F}[b]))\}$
10:     $\mathcal{L}_{CE_k}(\mathcal{Q}_{b,k}, \mathcal{Y}_{b,k}) \leftarrow$ Eq. 2, Cross-entropy loss
11:     $\theta_k^t \leftarrow \theta_k^t - \eta \nabla_{\theta_k^t} \mathcal{L}_{CE_k}$
12:    $\nabla\theta_k^t \leftarrow \theta_k^t - \theta_g^t$   ▷ **FAVD verified updates**
13:    **return** $\nabla\theta_k^t, \rho_k = (\tilde{\mu}_k, \tilde{\Sigma}_k)$
14: **Server execution** $(\nabla\theta_k^t, \rho_k)$:
15: Receive model updates from selected clients $\leftarrow \nabla\theta_k^t$
16: Perform model aggregation using FedAvg (**Eq. 1**)
17: Update the global model parameter: $\theta_g^{t+1}$
18: Update $\rho_g \leftarrow (\frac{1}{n}\sum_{k=1}^n \tilde{\mu}_k, \frac{1}{n}\sum_{k=1}^n \tilde{\Sigma}_k)$
19: Share $\theta_g^t, \rho_g = (\mu_g, \Sigma_g)$ to all the clients
20: Compute $\mathcal{A}_\mathcal{G} \leftarrow \text{Test}(\mathcal{G}_{\theta_g^{t+1}}, \mathcal{X}_{test})$
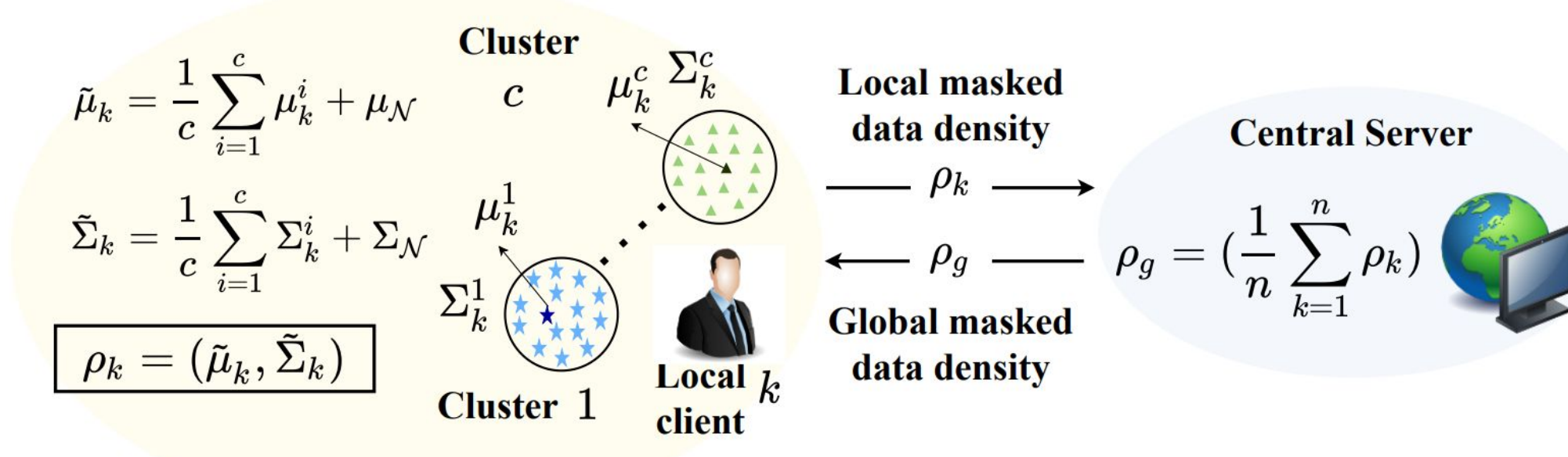21: **return** $\mathcal{A}_\mathcal{G}$



Figure 2. FAVD auditable data valuation process

**Definition 3.1** (ω-*bounded poisoned data.*) *Let $\mathcal{F}$ be the adversary's function class. An adversarial perturbation $\mathcal{U}$ is characterized by the mapping $\mathcal{U} := \mathcal{F} \times \mathcal{X}_k \times \mathbb{R} \to \Phi(\mathcal{X}_k)$. For $\omega > 0$, we define the $L_2$ norm ball as $\mathcal{B}_2(\mathcal{X}_k, \omega) := \{\Phi(\mathcal{X}_k) \in \mathbb{R}^d : \|\Phi(\mathcal{X}_k) - \mathcal{X}_k\| \leq \omega\} \bigcap \mathcal{X}_k$. We classify the adversarial perturbation $\mathcal{U}$ as ω-bounded if it adheres to the condition $\mathcal{U}(\mathcal{F}, \mathcal{X}_{k,i}, \mathcal{Y}_{k,i})_{i=1}^\nu \in \mathcal{B}(\mathcal{X}_{k,i}, \omega)_{i=1}^\nu$. Furthermore, for a given $\omega > 0$, we denote the worst-case adversarial perturbation $(\mathcal{U}^+)$ as*

$$\mathcal{U}^+ := \underset{\Phi(\mathcal{X}_k)\in\mathcal{B}(\mathcal{X}_k,\omega)}{\arg\max} \mathcal{L}_{CE}(f_\theta(\Phi(\mathcal{X}_{k,i}), \mathcal{Y}_{k,i})_{i=1}^\nu),$$

*where $\mathcal{L}_{CE}$ and $f_\theta$ represent the cross-entropy loss function and parameters of the local black-box model, respectively.*

**Table 1.** Global test accuracy (AG%) ↑ without data poisoning in uniform and non-IID FL settings, with and without FAVD data valuation. "No Val" refers to FL without data valuation.

| Dataset → | | C-xray [11] | Cam17 [3] | HAM10K [49] | CIFAR10 [24] | CIFAR100 [24] |
|---|---|---|---|---|---|---|
| **Shard ↓** | FL setting → Method ↓ | $n = 4$, $k = 4$ | $n = 5$, $k = 5$ | $n = 64$, $k = 64$ | $n = 100$, $k = 70$ | $n = 10^3$, $k = 500$ |
| **Uniform** | No Val | 91.37±1.54 | 91.82±1.21 | 73.61±1.59 | 85.59±0.76 | 78.45±0.81 |
| | **FAVD (ours)** | 93.75±1.20 | 94.62±1.15 | 75.48±0.42 | 88.42±0.21 | 81.19±0.45 |
| **non-IID** | No Val | 88.62±1.27 | 87.81±1.79 | 68.79±0.83 | 82.93±1.18 | 73.56±1.74 |
| | **FAVD (ours)** | 91.36±1.31 | 91.38±1.86 | 71.36±0.86 | 86.24±1.22 | 76.38±1.83 |

## Result Discussion

- **Higher accuracy in no-threat:** FAVD improves accuracy by 2-3% (uniform) and 3-4% (non-IID) by discarding misaligned data.
- **Resilient to poisoning attacks:** Limits accuracy drop to 3-7% (uniform) and 8-10% (in highly non-IID settings).
- **Superior data valuation:** Achieves highest MSDR (0.90) and BMR (0.86), surpassing other methods.
- **Better client contribution consistency (CCC):** Maintains high CCC (~1), ensuring verifiable contributions even under threats.

## Limitations and Future Work

- **Benign data exploitation:** Enhance feature alignment, explore encryption & secure multi party computation protocols.
- **Adaptive attack evasion:** Add behavioral analysis & temporal consistency checks.

## Summary and Conclusion

FAVD enhances auditability, verifiability, and privacy in FL by leveraging local and global density functions for precise data valuation. Extensive evaluations show superior performance over state-of-the-art methods across diverse datasets.

## References

K. Naveen Kumar, *et al.*, "The Impact of Adversarial Attacks on Federated Learning: A Survey," *IEEE TPAMI*, **2024**.

[1]{cs19m20p000001@, ckm@cse.}iith.ac.in, [2]rrjha@iitp.ac.in, [3]ravindrababu.t@gmail.com