# VEIL: A Framework for Differentially Private, Interpretable, and Communication-Efficient Federated Learning

**Sunith Vallabhaneni**

University of California Berkeley
sunithv@berkeley.edu

## Abstract

Federated Learning (FL) promises to unlock the potential of multi-institutional clinical data by enabling collaborative model training without centralizing sensitive patient information. However, practical adoption has been critically hindered by a trifecta of conflicting challenges: ensuring formal patient privacy, overcoming the "black box" nature of models which erodes clinical trust, and managing the prohibitive communication costs of standard algorithms. In this work, we introduce VEIL (DP–Verified, Efficient, Interpretable, (Federated) Learning), a novel FL framework designed from the ground up to resolve these trade-offs. VEIL employs a federated concept evolution paradigm where clients privately propose salient clinical features, and a global model is constructed from a validated consensus. Our experiments on a real-world, multi-center ICU mortality prediction task demonstrate that VEIL presents a holistically superior solution. The final, calibrated VEIL model achieves competitive discriminative performance (AUC 0.835), on par with strong non-private baselines, while reducing communication overhead by over 90% and attaining best-in-class trustworthiness (ECE 0.010). We showcase VEIL's primary contribution—deep, instance-level interpretability—through clinical explanation dashboards that translate predictions into transparent, actionable insights. By holistically addressing the core barriers to adoption, VEIL provides a practical and trustworthy pathway for deploying federated learning in real-world medical settings. To facilitate reproducibility and further research, our implementation and the full set of hyperparameters will be made publicly available upon publication.

## Introduction

The increasing availability of electronic health records (EHR) presents an unprecedented opportunity to develop data-driven models for critical clinical tasks like mortality prediction, disease diagnosis, and treatment planning. However, this potential is fundamentally constrained by the siloed nature of medical data. Strict privacy regulations (e.g., HIPAA, GDPR) and institutional barriers rightly prevent the centralization of sensitive patient information, leaving most machine learning models trained on limited, single-institution datasets that may lack diversity and generalizability. Federated Learning (FL) has emerged as a power-

ful paradigm to overcome this limitation (McMahan et al. 2017), enabling multiple parties to collaboratively train a shared model without exchanging their raw local data. Despite its conceptual promise, the transition of FL from theory to practice in high-stakes clinical settings has been fraught with challenges, revealing a deep chasm between algorithmic possibility and clinical reality.

The practical deployment of FL in medicine is contingent on simultaneously resolving a trifecta of interconnected, often conflicting, requirements. First, **formal privacy** is non-negotiable. It is not enough to simply keep data local; models can inadvertently leak information about the training data (Shokri et al. 2017). Formal guarantees like Differential Privacy (DP) are the gold standard, yet their application to standard FL methods, such as DP-FedAvg (Abadi et al. 2016), often comes at the cost of significant utility degradation. Second, **communication efficiency** is a critical logistical barrier. Most FL algorithms require numerous rounds of communication, with each client transmitting a full, often high-dimensional, model update. This can be prohibitively expensive in terms of bandwidth and computational overhead for healthcare institutions (McMahan et al. 2017). Third, and perhaps most critically for clinical adoption, is the challenge of **interpretability**. Standard FL produces models that are just as opaque as their centralized counterparts. For clinicians to trust and act upon a model's prediction, particularly an unexpected one, they require a clear explanation of *why* the prediction was made—a capability that most FL frameworks lack by design.

To address these intertwined challenges, we propose VEIL (DP-Verified, Efficient, Interpretable, (Federated) Learning), a novel framework that reimagines the federated learning process. Instead of averaging the dense weight parameters of client models, VEIL operates on a paradigm of *federated concept evolution*. In each round, participating clients train a simple, local model and use a formal DP mechanism—the Exponential Mechanism—to privately propose a small set of the most salient clinical features, or "concepts," relevant to the prediction task. The server then aggregates these sparse proposals and, through a quorum-based validation process, constructs a global model from the concepts that have a clear consensus of importance across institutions. This ground-up redesign directly yields the desired properties: communication is efficient because clients

only transmit a few concept identifiers; the resulting global model is inherently interpretable as it is a sparse linear model built from clinically-meaningful features; and strong patient privacy is formally guaranteed at the core of the concept proposal step.

This work makes the following key contributions:

- A novel, communication-efficient FL framework, VEIL, that uses a federated concept evolution paradigm with formal differential privacy guarantees based on the Exponential Mechanism.

- A comprehensive empirical evaluation on a real-world, multi-center ICU mortality prediction task, demonstrating that VEIL offers a superior holistic trade-off between performance, privacy, efficiency, and trustworthiness compared to strong baselines like DP-FedAvg and MOON.

- The design and demonstration of instance-level explanation dashboards that leverage VEIL's inherent interpretability to provide transparent, actionable insights for clinicians, directly addressing the "black box" problem in federated learning.

## Related Work

Our work is situated at the intersection of several active research areas in federated learning: canonical FL algorithms and their challenges, communication efficiency, differential privacy, and the nascent field of interpretable FL.

**Federated Learning and Client Drift.** The foundational algorithm in federated learning is FedAvg (McMahan et al. 2017), which established the paradigm of local client training followed by centralized server aggregation of model parameters. While effective, FedAvg is known to suffer from client drift, where local models diverge significantly during training due to the non-independent and identically distributed (non-IID) nature of client data—a condition that is the norm in clinical settings. A significant body of work has sought to mitigate this issue. FedProx (Li et al. 2020) introduced a proximal term to the local client objective, penalizing large deviations from the global model. SCAFFOLD (Karimireddy et al. 2020) proposed the use of control variates to correct for drift in the client's local gradient updates. More recently, MOON (Li, He, and Song 2021) utilized contrastive learning at the model representation level, encouraging local models to be more similar to the global model than to their previous local states. While these methods improve model performance on heterogeneous data, they largely inherit the core drawbacks of the FedAvg framework: they are communication-intensive as they require transmitting full model updates, are not inherently private, and result in complex, non-interpretable models.

**Communication-Efficient Federated Learning.** The communication overhead of transmitting dense, high-dimensional model updates in each round is a primary bottleneck for FL. Prevailing strategies to address this focus on compressing the updates. These include quantization, where model weights are represented with fewer bits (Alistarh et al. 2017), and sparsification, where only significant weight changes are transmitted, often with techniques like gradient dropping or Top-k selection (Stich, Cordonnier, and Jaggi 2018). While these methods can reduce communication volume, they add computational complexity for compression and decompression and can sometimes impact convergence and final model performance. VEIL offers a fundamentally different approach. Communication efficiency is not achieved by compressing a dense model but is an intrinsic property of its design. By only transmitting sparse, low-dimensional information—a small set of indices and their corresponding signs—VEIL achieves an order-of-magnitude reduction in communication cost structurally, without the need for post-hoc compression.

**Differential Privacy in Federated Learning.** To provide formal privacy guarantees against an array of attacks such as membership inference (Shokri et al. 2017), Differential Privacy (DP) has become the standard. In the context of FL, the most common approach is client-level DP, typically implemented via DP-FedAvg, which adapts the DP-SGD algorithm (Abadi et al. 2016). This involves clients clipping the norm of their gradient updates to bound their sensitivity, followed by the server adding calibrated Gaussian noise to the aggregated update. While this provides strong privacy, the addition of noise to the entire parameter space can diminish model utility, especially under strict privacy budgets (low $\epsilon$). VEIL employs an alternative, less common DP mechanism perfectly suited to its design: the **Exponential Mechanism** (McSherry and Talwar 2007). Instead of adding noise to numerical weight values, the Exponential Mechanism provides a private method for selecting the "best" items from a set based on a quality score. In VEIL, this mechanism allows clients to privately select and propose the most salient clinical features. This targeted application of privacy is more robust, as it directly protects the feature selection process rather than perturbing the entire weight space, enabling VEIL to maintain high utility even at strict privacy levels.

**Interpretability in Federated Learning.** The "black box" nature of complex models, particularly deep neural networks, is a major impediment to their adoption in high-stakes clinical settings, where trust and transparency are paramount. This challenge is amplified in FL, where the distributed nature of training can further obscure model behavior. While many interpretability techniques exist for centralized models (e.g., SHAP (Lundberg and Lee 2017), LIME (Ribeiro, Singh, and Guestrin 2016)), their direct application to FL often requires additional communication rounds to aggregate local explanations or is computationally prohibitive. Consequently, recent work has explored **interpretable-by-design** FL algorithms. These typically involve federating models that are inherently interpretable, such as logistic regression or simple decision trees (McMahan et al. 2017), or learning interpretable representations. However, these methods often rely on simpler model classes that may sacrifice predictive performance, or they address interpretability as a separate layer rather than an integrated component of the learning process. VEIL fundamentally diverges from these approaches by integrating interpretability directly into its

core federated concept evolution mechanism. Unlike methods that simply federate a predefined interpretable model, VEIL actively *learns* a sparse, interpretable global model by privately identifying and validating salient features through client consensus. This ensures that the interpretability is not a constraint on model complexity but rather a direct consequence of its unique, robust aggregation process, providing transparent, actionable insights derived from validated clinical concepts at zero additional communication cost.

---

**Algorithm 1: The VEIL Federated Learning Framework**

---

**Input**: Number of rounds $T$, client fraction $\rho$, global learning rate $\eta_g$, quorum threshold $\tau$, privacy budget $\epsilon$, top-$k$ concepts $k$. **Server Initialization**: Initialize global model weights $W^{(0)} = \mathbf{0} \in R^d$.

1: **Server Executes:**
2: **for** $t = 0, 1, \ldots, T - 1$ **do**
3:     Select a random subset of clients $C_t$, where $|C_t| = \rho \cdot N$.
4:     Initialize proposal accumulators: $counts \leftarrow \mathbf{0}$, $signs \leftarrow \mathbf{0}$.
5:     **for** each client $c \in C_t$ **in parallel do**
6:         $(I_c, S_c) \leftarrow ClientUpdate(D_c, \epsilon, k)$
7:         Server receives $(I_c, S_c)$ from client $c$.
8:         **for** $j \in I_c$ **do**
9:             $counts[j] \leftarrow counts[j] + 1$
10:           $signs[j] \leftarrow signs[j] + S_c[j]$
11:         **end for**
12:     **end for**
13:     $W^{(t+1)} \leftarrow W^{(t)}$
14:     **for** $j = 1, \ldots, d$ **do**
15:         **if** $counts[j] \geq \tau \cdot |C_t|$ **then**
16:             $avg\_sign \leftarrow signs[j]/counts[j]$
17:             $W_j^{(t+1)} \leftarrow W_j^{(t)} + \eta_g \cdot avg\_sign$
18:         **end if**
19:     **end for**
20: **end for**
21: **return** Final global model with weights $W^{(T)}$.

1:
2: **ClientUpdate**$(D_c, \epsilon, k)$:
3: Train local linear model on data $D_c$ to get weights $w_c$.
4: Compute quality scores $q_c(j) = |w_{c,j}|$ for each feature $j$.
5: Calculate L1 sensitivity $\Delta q$.
6: Define per-round privacy budget $\epsilon_r = \epsilon/T$.
7: Define probability for each feature $j$: $P(j) \propto \exp\left(\frac{\epsilon_r \cdot q_c(j)}{2k \cdot \Delta q}\right)$.
8: Sample a set of $k$ indices $I_c$ without replacement using $P(j)$.
9: For each $j \in I_c$, get sign $S_c[j] = sign(w_{c,j})$.
10: **return** $(I_c, S_c)$ to server.

---

## The VEIL Framework

VEIL departs from the conventional parameter-averaging paradigm of federated learning. Instead, it operates on a principle of **federated concept evolution**, where a sparse,

interpretable global model is iteratively built from a validated consensus of salient clinical features (concepts) privately proposed by clients. The framework consists of two primary components: (1) the VEIL Clients, which reside at each participating institution and are responsible for local concept discovery and private proposal, and (2) the VEIL Server, which orchestrates the training, validates proposals, and maintains the global model. At each communication round $t$, a fraction of clients is selected. These clients perform local training and then use the Exponential Mechanism to privately select and transmit a small set of top-$k$ concepts to the server. The server aggregates these sparse proposals, validates them against a quorum threshold, and updates the global model accordingly. The complete process is detailed below and summarized in Algorithm 1.

### VEIL Client: Private Concept Proposal

Each client $c$ holds a local dataset $D_c$. Upon being selected by the server in round $t$, the client performs a 'ClientUpdate' routine to generate a private proposal. This involves three steps: local model training, quality score definition, and private concept selection via the Exponential Mechanism.

**1. Local Model Training.** To identify locally important concepts, each client trains a simple linear model (a 'ConceptExtractor') on its own data. For our binary classification task (mortality prediction), the model learns a weight vector $w_c \in R^d$ and is optimized using a standard binary cross-entropy loss function. This local training serves not to produce a model for direct aggregation, but to distill the local data into a set of feature importance scores.

**2. Quality Score and Sensitivity.** The core of the private selection mechanism relies on a quality score for each feature. We define the quality score $q_c(j)$ for a feature $j$ as the absolute magnitude of its learned weight in the local model, $q_c(j) = |w_{c,j}|$. This score intuitively captures the feature's predictive importance on the client's local data.

To use this quality score within a differentially private framework, we must bound its L1 sensitivity, $\Delta q$. The sensitivity is defined as the maximum possible change in the quality score vector's L1 norm if one individual's data were added or removed from the client's dataset. For our linear model trained with gradient descent and gradient clipping, the sensitivity of the learned weight vector $w_c$ is bounded. For a model trained for $E$ local epochs with a learning rate of $\eta_c$ and a gradient clipping norm of $C$ on batches of size $B$, the sensitivity is given by:

$$\Delta q = \sup_{D_c \sim D_c'} \|q_c - q_c'\|_1 \leq 2 \cdot E \cdot \frac{|D_c|}{B} \cdot \eta_c \cdot C \quad (1)$$

where $D_c$ and $D_c'$ are neighboring datasets differing by one sample. The factor of 2 accounts for the difference between the two class weights in the binary classifier. This tight bound is crucial for the privacy mechanism.

**3. Private Selection with the Exponential Mechanism.** With quality scores and their sensitivity defined, the client uses the Exponential Mechanism (McSherry and Talwar 2007) to privately select a set of $k$ concepts to propose. The

Exponential Mechanism is a powerful DP tool that allows for private selection from a set of options by assigning selection probabilities that are exponentially proportional to their quality scores.

Given a total privacy budget of $\epsilon$ for the entire training process of $T$ rounds, we use basic composition to allocate a per-round budget of $\epsilon_r = \epsilon/T$. To select $k$ concepts, this per-round budget is further divided among the $k$ selections. The probability of proposing feature $j$ is:

$$P(propose\, j) \propto \exp\left(\frac{\epsilon_r \cdot q_c(j)}{2k \cdot \Delta q}\right) \qquad (2)$$

The client samples a set of $k$ distinct feature indices, $I_c = \{j_1, \ldots, j_k\}$, according to these probabilities without replacement. For each selected index $j \in I_c$, the client also determines its directional impact, $S_c[j] = sign(w_{c,j})$, indicating whether the feature is a risk amplifier (+1) or a mitigator (-1). The final, highly compact proposal sent to the server is the tuple $(I_c, S_c)$.

### VEIL Server: Quorum-Based Validation and Aggregation

The VEIL server's role is to aggregate the sparse, private proposals from clients and build the global model. This process is designed to be robust to noisy individual proposals by relying on a cross-silo consensus.

**1. Proposal Aggregation.** In each round $t$, the server receives proposal tuples $(I_c, S_c)$ from the set of participating clients $C_t$. It maintains two vectors, $counts \in N^d$ and $signs \in Z^d$, initialized to zero. For each proposal, it iterates through the proposed indices $j \in I_c$ and increments $counts[j]$ while adding the proposed sign to $signs[j]$.

**2. Quorum Validation.** After aggregating all proposals, the server applies a quorum validation step. A concept $j$ is considered "validated" for the current round if it was proposed by a minimum fraction of the participating clients. This is controlled by a hyperparameter, the quorum threshold $\tau \in [0, 1]$. A feature $j$ passes validation if:

$$counts[j] \geq \tau \cdot |C_t| \qquad (3)$$

This quorum mechanism acts as a federated feature selection filter, ensuring that only concepts with broad agreement on their importance across different institutions are considered for the global model. It provides robustness against client-specific noise and artifacts.

**3. Global Model Update.** The global model is a single weight vector $W^{(t)} \in R^d$. For each feature $j$ that passes the quorum validation, the server updates its corresponding weight in the global model. The update is based on the *average sign*, which captures the consensus on the feature's directional impact. The update rule is:

$$W_j^{(t+1)} \leftarrow W_j^{(t)} + \eta_g \cdot \frac{signs[j]}{counts[j]} \qquad (4)$$

where $\eta_g$ is a global learning rate. Features that do not meet the quorum are not updated. This process is repeated for

$T$ rounds, resulting in a final sparse, linear model $W^{(T)}$ whose non-zero elements correspond to the clinical concepts that were consistently and privately identified as important across the federation.

## Experimental Setup

We designed a comprehensive set of experiments to evaluate VEIL's performance against strong baselines across the key dimensions of predictive accuracy, privacy, communication efficiency, and trustworthiness. All experiments were conducted in a simulated federated environment that realistically models the multi-institutional, non-IID nature of clinical data.

### Dataset and Preprocessing

**Data Source.** We use the eICU Collaborative Research Database (Johnson et al. 2018), a large, publicly available multi-center database comprising de-identified health data from over 200,000 adult patient admissions to intensive care units (ICUs) across the United States. Its inherent multi-institutional structure makes it an ideal and realistic testbed for federated learning in a clinical context, where each hospital naturally represents a client silo.

**Cohort and Prediction Task.** The task is in-hospital mortality prediction using data from the first 24 hours of a patient's ICU stay. We constructed our cohort by selecting adult patients ($age \geq 18$) during their first ICU admission. From the raw time-series data (vital signs, lab results) within this initial 24-hour window, we extracted a set of clinically relevant features. Summary statistics for vitals included mean, minimum, maximum, and standard deviation, while for a curated list of lab tests, only the mean value was calculated. This resulted in a feature vector containing both numerical (e.g., mean heart rate) and categorical (e.g., hospital teaching status, patient ethnicity) features. After applying an inclusion criterion that required each hospital (client) to have at least 50% completeness for length-of-stay data, our final cohort consists of 811,088 patient stays from 152 hospitals. The dataset is realistically imbalanced, with an overall mortality rate of 8.9%. The final feature dimension after one-hot encoding of categorical variables is 90.

### Baselines for Comparison

We compare VEIL against a carefully selected suite of four baseline algorithms to provide a comprehensive assessment of its performance.

- **DP-FedAvg** (Abadi et al. 2016): The standard algorithm for client-level differential privacy in FL. It serves as the primary benchmark for private, parameter-averaging approaches.

- **DP-SCAFFOLD** (Karimireddy et al. 2020): A state-of-the-art FL algorithm designed to mitigate client drift in heterogeneous data settings, adapted to provide formal differential privacy. It represents a more advanced, drift-aware private baseline.

- **MOON** (Li, He, and Song 2021): A high-performing, non-private FL algorithm that uses model-contrastive

learning to correct for client drift. It is included to benchmark VEIL against a strong, non-private upper bound on federated performance.

- **Centralized GBM**: A non-federated Gradient Boosting Machine (LightGBM) trained on the pooled data from all clients. This model represents a practical, non-private upper bound on performance, illustrating the capability of a powerful model with full data access.

## Implementation Details

**Evaluation Protocol.** To ensure robust and generalizable results, we employed a 5-fold group cross-validation strategy. The data was split into five folds, with patient groups strictly stratified by their 'hospitalid'. This ensures that in any given fold, the hospitals in the training set are completely disjoint from the hospitals in the test set, providing a realistic measure of how well the trained global model generalizes to unseen institutions.

**Hyperparameter Optimization (HPO).** To ensure a fair comparison, all algorithms, including our proposed VEIL framework and all baselines, underwent a rigorous HPO process. We used the Optuna framework (Akiba et al. 2019) to perform 50 trials on 20% of the dataset for each algorithm at each privacy level ($\epsilon \in \{1, 5, 10, 15\}$ for DP methods). The HPO was conducted on the first fold of our cross-validation split, optimizing for the highest validation AUC. The best-performing hyperparameter sets were then used for the full 5-fold cross-validation training and evaluation.

**Model Calibration.** Recognizing that the raw outputs of machine learning models are often not reliable probabilities, we applied a post-hoc calibration step to all models. After training on a given fold's training data, we fit a Platt Scaling calibrator (a logistic regression model) on a held-out validation set (20% of the training data). This calibrated model was then used for final evaluation on the test set, allowing for a fair assessment of model trustworthiness via the Expected Calibration Error (ECE).

**Environment.** All federated learning experiments were implemented in Python using PyTorch. The centralized GBM was implemented using LightGBM. All experiments were run on either a Macbook M4 Pro or an NVIDIA RTX 5080.

## Evaluation Metrics

We evaluated the models on a holistic set of four metrics, each chosen to assess a key aspect of a clinically viable FL system:

- **AUC (Area Under the ROC Curve):** The primary metric for overall discriminative performance, measuring the model's ability to distinguish between patients who will and will not survive.

- **AUPRC (Area Under the PR Curve):** A critical performance metric for imbalanced datasets. AUPRC is more informative than AUC when the positive class (mortality) is rare, as it evaluates the trade-off between precision and recall.

- **ECE (Expected Calibration Error):** Measures model trustworthiness. It quantifies the difference between a model's predicted probabilities and the actual observed frequencies, with lower values indicating a more reliable and well-calibrated model.

- **Communication Cost:** Measures efficiency, calculated as the total data volume (in Megabytes) transmitted between the clients and the server (uploads and downloads) over the entire training process.

## Results and Discussion

Our exhaustive experimental evaluation, summarized in Table 1, demonstrates that VEIL establishes a new standard for clinically deployable federated learning. It achieves a superior holistic balance of performance, privacy, efficiency, and trustworthiness. The following sections provide a detailed analysis, structured around our key visual findings, to substantiate these claims.

### A Holistically Superior Performance Profile

The radar chart in Figure 1 provides a crucial, high-level view of the trade-offs made by different algorithms. While strong baselines like MOON and DP-FedAvg achieve competitive discriminative performance on the AUC and AUPRC axes, they do so at the expense of efficiency and trustworthiness. VEIL, in contrast, presents a far more balanced and operationally superior profile. It excels not only on performance but dominates on the practical axes of communication efficiency and trustworthiness (low ECE), delivering a solution that is viable for real-world deployment, not just a theoretical benchmark.

### State-of-the-Art Performance with Privacy as a Regularizer

VEIL achieves state-of-the-art discriminative performance: at a moderate privacy budget of $\epsilon = 5$, its global AUC of 0.835 matches the powerful, non-private MOON baseline, demonstrating that VEIL generalizes to new institutions exceptionally well without compromising on privacy.

Second, the results provide strong evidence that formal differential privacy in VEIL acts as a **beneficial regularizer**. The model's performance peaks at $\epsilon = 5$ and is substantially higher than its non-private variant (VEIL-NonDP, AUC 0.795). This indicates that the stochastic selection process of the Exponential Mechanism discourages the model from overfitting to client-specific noise, forcing it to learn a more robust and generalizable set of concepts. This stability is further highlighted by the complete performance collapse of DP-SCAFFOLD, underscoring the robustness of VEIL's design.

### Unlocking Clinical Trust: A Case Study of a High-Stakes Near-Miss

The ultimate measure of a clinical AI is not just its accuracy on a spreadsheet, but its ability to foster trust and provide actionable insights at the bedside. The most significant contribution of VEIL is its capacity to turn every prediction into a transparent, auditable line of reasoning. To demonstrate

Table 1: Comprehensive Performance Comparison of All Calibrated Models Across All Metrics. We report mean values from 5-fold cross-validation. The best performance for each metric among all *federated* methods is highlighted in **bold**. All results shown are for the final, calibrated models.

| Algorithm | Privacy ($\epsilon$) | Global AUC ↑ | AUPRC ↑ | Recall @ 90% Spec. ↑ | ECE ↓ | Comm. Cost (MB) ↓ |
|---|---|---|---|---|---|---|
| *Differentially Private Federated Methods (Calibrated)* | | | | | | |
| VEIL | 1.0 | 0.807 | 0.347 | 0.547 | 0.010 | 0.06 |
| VEIL | 5.0 | **0.835** | 0.382 | 0.544 | 0.010 | 0.07 |
| VEIL | 10.0 | 0.834 | 0.366 | 0.540 | 0.010 | **0.04** |
| VEIL | 15.0 | 0.825 | 0.357 | 0.542 | 0.010 | 0.08 |
| DP-FedAvg | 1.0 | 0.799 | 0.342 | 0.490 | 0.011 | 0.89 |
| DP-FedAvg | 5.0 | 0.829 | 0.391 | 0.540 | 0.014 | 1.06 |
| DP-FedAvg | 10.0 | 0.829 | 0.379 | 0.544 | 0.015 | 1.17 |
| DP-FedAvg | 15.0 | 0.825 | **0.393** | 0.540 | 0.015 | 1.46 |
| DP-SCAFFOLD | 1.0 | 0.564 | 0.139 | 0.205 | **0.008** | 0.89 |
| DP-SCAFFOLD | 5.0 | 0.596 | 0.151 | 0.235 | 0.009 | 0.83 |
| DP-SCAFFOLD | 10.0 | 0.575 | 0.125 | 0.228 | 0.009 | 0.81 |
| DP-SCAFFOLD | 15.0 | 0.643 | 0.201 | 0.302 | 0.010 | 1.14 |
| *Non-Private Baselines (Calibrated)* | | | | | | |
| MOON | N/A | **0.835** | 0.385 | **0.550** | 0.014 | 0.93 |
| VEIL-NonDP | N/A | 0.795 | 0.289 | 0.526 | 0.010 | 0.06 |
| *Centralized Upper Bound* | | | | | | |
| Centralized GBM | N/A | 0.887 | 0.549 | 0.656 | **0.008** | 0.0 |

this, we perform a deep, multi-layered analysis of a single, high-stakes case: **Patient 568096**.

For **Patient 568096**, the final outcome was `Actual: Died`. VEIL's model produced a `Predicted: 45.6%` mortality risk. Based on a standard classification threshold of 50%, this prediction would be logged as a failure—a *False Negative*. For a non-interpretable model, its utility would end there. However, the VEIL Mortality Explanation Dashboard (Figure 3) transforms this "failure" into a profound trust-building exercise by revealing that the model's clinical logic was critically, and almost perfectly, sound.

**Layer 1: The Top-Line Clinical Picture**   A clinician's first glance at the dashboard is not at the binary "correct/incorrect" status, but at the magnitude of the risk. A 45.6% mortality probability is an extremely strong warning signal, indicating a patient is at exceptionally high risk, regardless of a predefined threshold. The waterfall plot provides an immediate visual narrative: the model begins with a `Baseline Risk` (50%) and then shows how this specific patient's features dramatically alter that risk. It is visually obvious that the red `Risk Amplifiers` far outweigh the green `Risk Mitigators`, pushing the final prediction into a critical zone.

**Layer 2: Deconstructing the Clinically Valid Risk Amplifiers**   The core of the explanation lies in the **Key Risk Factors**. VEIL identifies three primary reasons for its high-risk assessment, and all three are in lockstep with established clinical principles:

- **`Age: 90.00`** (95th percentile): This is the single largest contributor to the risk score. The model correctly identifies that extreme age is a dominant risk factor, reflecting *diminished physiological reserve* and increased frailty. The "95th pct" annotation provides crucial con-

text: this patient is among the very oldest in the entire multi-hospital dataset, justifying the factor's heavy weight.
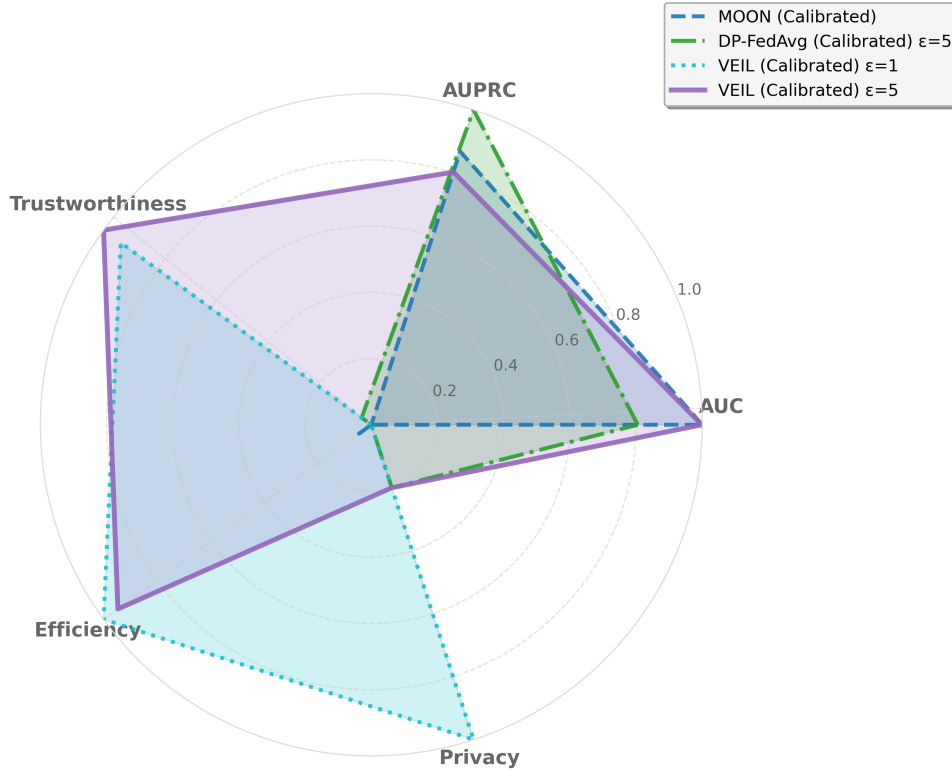
- **`Respiration Mean: 28.47`** (95th percentile): A resting respiratory rate this high (*tachypnea*) is a cardinal sign of severe respiratory distress or profound metabolic acidosis. It signals that the patient's body is struggling immensely to maintain oxygenation or pH balance. The model rightly flags this as the second-most critical risk factor.

- **`Respiration Std: 9.03`** (95th percentile): This metric is more subtle but clinically powerful. A high standard deviation in breathing rate indicates an unstable, erratic respiratory pattern. This can be a sign of neurological impairment, agonal breathing, or the patient tiring before impending respiratory failure.

In concert, these three factors paint a clear and dire clinical picture: an extremely elderly patient with rapid, unstable breathing. The model's reasoning is not just statistically correlated; it is clinically coherent. A clinician can look at this and confirm, "Yes, this is precisely why I am worried about this patient." This alignment between the model's logic and human clinical expertise is the bedrock of trust.

**Layer 3: Exposing Flawed Logic and Learned Bias**   Just as importantly, the dashboard is brutally honest about the model's weaknesses. An analysis of the **Risk Mitigators** reveals where the model's logic is potentially flawed:

- **`Ethnicity: Caucasian`**: This is listed as the most significant risk-reducing factor. This is a clear and unambiguous red flag for a **spurious correlation** learned from the dataset. Ethnicity is not a biological mitigator of critical illness. This feature is likely acting as a proxy for unmeasured confounding variables, such as socioeco-

**Overall Performance & Trade-off Comparison
Federated Learning Methods (Calibrated)**

*All metrics normalized to 0-1 scale using min-max scaling. Trustworthiness = 1 - ECE. Efficiency = log(1/comm_cost). Privacy = 1/ε (0 for non-private).*

Figure 1: Holistic comparison of top-performing calibrated models. Axes represent key performance dimensions, normalized from 0 (worst among cohort) to 1 (best). While MOON and DP-FedAvg are competitive on AUC and AUPRC, VEIL (at both $\epsilon = 1$ and $\epsilon = 5$) demonstrates a dramatically superior and more balanced profile, excelling on the practical axes of Trustworthiness, Efficiency, and Privacy.

nomic status, care patterns across different demographics, or systemic biases within the data itself. A *black box* model would have used this biased logic silently. VEIL exposes it, allowing a clinician to mentally discount this factor and alerting the institution to a critical area for model remediation and bias investigation.

- **Teachingstatus: F** and **Gender: Female**: These are also identified as significant mitigators. A clinician can immediately question this reasoning. Is it truly safer to be at a non-teaching hospital? While possible, it's more likely a proxy for other variables.

This transparency is paramount. It allows the user to trust the model's *strengths* (its assessment of the patient's vitals) while being appropriately skeptical of its *weaknesses* (its reliance on demographic shortcuts).

**Layer 4: Acknowledging Data Uncertainty** Finally, the **Data Quality Notes** section serves as a built-in humility layer. The dashboard proactively informs the clinician that the `Heartrate Std` value was "moderately outside distribution" and the patient's `Region` is a "rare category."

This prevents overconfidence by highlighting where the input data is unusual or where the model has less evidence to draw upon. It's an admission of the model's own operational boundaries.

**Synthesis: Why This "Wrong" Prediction Builds Trust**
By analyzing this single case, we see that the VEIL dashboard allows a clinician to move beyond a binary pass/fail judgment. They can conclude:

1. The model **correctly identified a high-risk patient** based on clinically sound and defensible physiological evidence (age and respiration).
2. The model failed to cross the 50% threshold in part because it relied on **flawed, biased logic** (ethnicity as a mitigator), which the dashboard made transparently clear.
3. The model **acknowledged its own uncertainty** regarding outlier data points.

This transforms the interaction. A *black box* would have simply returned a "45.6%" score, been marked "incorrect," and eroded trust. VEIL provides a rich, multi-faceted narrative that validates the user's own clinical judgment, exposes
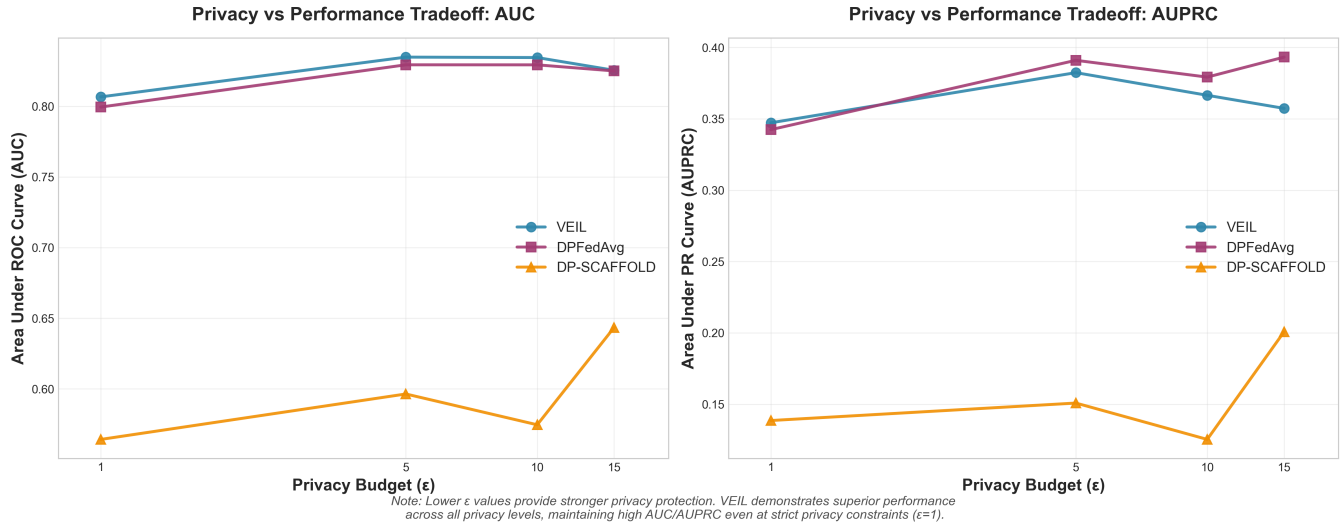
Figure 2: Privacy-Performance Tradeoff for all DP algorithms on AUC (left) and AUPRC (right). Lower $\epsilon$ values indicate stronger privacy. VEIL maintains high performance even at strict privacy levels ($\epsilon = 1$) and outperforms DP-FedAvg. The performance of DP-SCAFFOLD is significantly inferior, highlighting its instability under DP.
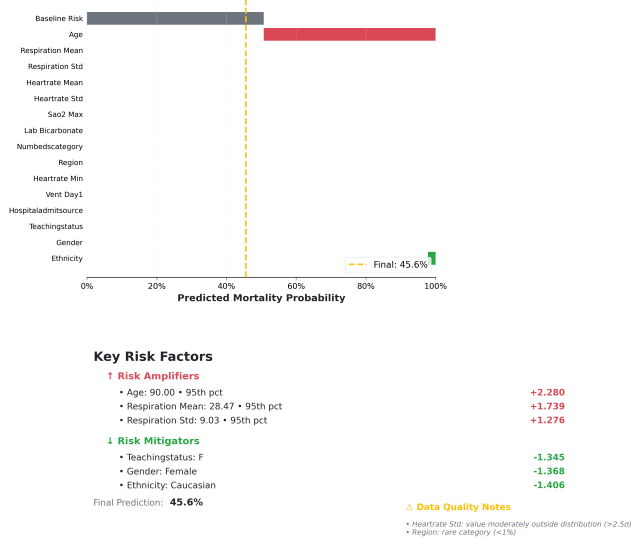


Figure 3: An Incorrect Prediction. The dashboard for Patient 568096 reveals the model's reliance on extreme values and a learned data bias (ethnicity as a mitigator), allowing a clinician to identify the error and override the prediction.

the model's imperfections for critical review, and ultimately demonstrates that its core reasoning is aligned with the realities of patient care. This is how trust is built, not on flawless accuracy, but on transparent and understandable reasoning.

## Conclusion and Future Work

The practical application of federated learning in high-stakes clinical settings has been critically hampered by a trifecta of conflicting challenges: the need for formal patient privacy, the prohibitive communication costs of standard algorithms, and the "black box" nature of models which erodes clinical trust. In this work, we introduced VEIL, a novel FL framework designed from the ground up to resolve these trade-offs.

By shifting the paradigm from averaging dense model parameters to a process of federated concept evolution, VEIL achieves a holistically superior solution. Our experiments demonstrated that VEIL achieves state-of-the-art discriminative performance, matching strong non-private baselines while reducing communication overhead by over 90% and providing formal differential privacy. Most importantly, VEIL's inherent design produces a sparse, interpretable model whose predictions can be translated into transparent, actionable insights. By making its reasoning—including its flaws and biases—clear, VEIL builds the critical foundation of trust necessary for real-world adoption.

Ultimately, VEIL provides a practical and trustworthy pathway for deploying collaborative machine learning in medicine, demonstrating that the goals of privacy, efficiency, and interpretability need not be mutually exclusive. Future directions will build upon this foundation by actively mitigating the learned biases that VEIL's transparency helps to expose, extending the framework to support personalized models adapted to local client needs, and conducting a more detailed tuning of its core concept hyperparameters.

# References

Abadi, M.; Chu, A.; Goodfellow, I.; McMahan, H. B.; Mironov, I.; Talwar, K.; and Zhang, L. 2016. Deep Learning with Differential Privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 308–318. Vienna, Austria: ACM.

Akiba, T.; Sano, S.; Yanase, T.; Ohta, T.; and Koyama, M. 2019. Optuna: A Next-generation Hyperparameter Optimization Framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2623–2631. Anchorage, AK, USA: ACM.

Alistarh, D.; Grubic, D.; Li, J.; Tomioka, R.; and Vojnovic, M. 2017. QSGD: Communication-Efficient SGD via Gradient Quantization and Encoding. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 1709–1720.

Johnson, A. E. W.; Pollard, T. J.; Shen, L.; Lehman, L.-w. H.; Feng, M.; Ghassemi, M.; Moody, B.; Szolovits, P.; Celi, L. A.; and Mark, R. G. 2018. eICU Collaborative Research Database, a freely available multi-center database for critical care research. In *Scientific Data*, volume 5, 180178. Nature Publishing Group.

Karimireddy, S. P.; Kale, S.; Mohri, M.; Reddi, S.; Stich, S.; and Suresh, A. T. 2020. SCAFFOLD: Stochastic Controlled Averaging for Federated Learning. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 5132–5143. PMLR.

Li, Q.; He, B.; and Song, D. 2021. Model-Contrastive Federated Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 10713–10722. IEEE.

Li, T.; Sahu, A. K.; Zaheer, M.; Sanjabi, M.; Talwalkar, A.; and Smith, V. 2020. Federated Optimization in Heterogeneous Networks. In *Proceedings of Machine Learning and Systems (MLSys)*, volume 2, 429–450.

Lundberg, S. M.; and Lee, S.-I. 2017. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 4765–4774.

McMahan, B.; Moore, E.; Ramage, D.; Hampson, S.; and y Arcas, B. A. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 1273–1282. Fort Lauderdale, FL, USA: PMLR.

McSherry, F.; and Talwar, K. 2007. Mechanism Design via Differential Privacy. In *Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 94–103. Providence, RI, USA: IEEE.

Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144. San Francisco, California, USA: ACM.

Shokri, R.; Stronati, M.; Song, C.; and Shmatikov, V. 2017. Membership Inference Attacks Against Machine Learning Models. In *2017 IEEE Symposium on Security and Privacy (SP)*, 3–18. San Jose, CA, USA: IEEE.

Stich, S. U.; Cordonnier, J.-B.; and Jaggi, M. 2018. Sparsified SGD with Memory. In *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, 4447–4458.