



AL- Rafidain University

PISSN: (1681-6870); EISSN: (2790-2293)

**Journal of AL-Rafidain  
University College for Sciences**

Available online at: <https://www.jrucs.iq>

**JRUCS**

Journal of AL-Rafidain  
University College  
for Sciences

## Privacy-Preserving and Fair Federated Data Mining for Early Prediction of Health Crises: An Empirical Study on Baghdad Hospitals

Mayyadah J. Gailan 

[mayadajabbar@uomustansiriyah.edu.iq](mailto:mayadajabbar@uomustansiriyah.edu.iq)

Department of Hotel Studies, College of Tourism Sciences, Al-Mustansiriyah University, Baghdad, Iraq

### Article Information

#### Article History:

Received: April, 20, 2026

Accepted: May, 17, 2026

Available Online: June, 30, 2026

#### Keywords:

Federated Learning, Data Mining, Health Crisis Prediction, Algorithmic Fairness, Differential Privacy, Electronic Health Records, Baghdad Hospitals, Edge Computing, Equalized Odds, FedFairMine..

### Abstract

Early health crises, such as sepsis, cardiac decompensation, and respiratory failure, kill patients who might otherwise survive, and the problem is worse in resource-limited settings like Iraq, where data infrastructure is fragile, and populations are routinely underrepresented in predictive models. The study built FedFairMine to address this directly. It is a federated data mining framework that keeps raw patient records on local hospital servers while still training a shared predictive model. Three design choices drive it: gradient-based feature extraction runs entirely on-device, so no raw data leaves the hospital; a fairness-weighted aggregation step explicitly corrects for performance gaps across age, gender, and socioeconomic groups; and a differential privacy layer tuned specifically for Iraqi electronic health records. The study tested the system across six Baghdad teaching hospitals — Al-Kindi, Ghazi Al-Hariri, Ibn Sina, Medical City, Al-Yarmouk, and Al-Kadhimiya — using 47,312 de-identified patient episodes from 2021 to 2024. FedFairMine reached an F1-score of 0.912, beating FedAvg, FedProx, and Scaffold by 8.3%, 5.1%, and 3.7%. Demographic disparity dropped to an equalized odds difference of 0.021, a 76% improvement over centralized baselines. The results show that predictive accuracy, demographic fairness, and patient privacy are not competing goals; they can be achieved together.

#### Correspondence:

Mayyadah J. Gailan

[mayadajabbar@uomustansiriyah.edu.iq](mailto:mayadajabbar@uomustansiriyah.edu.iq)

DOI: <https://doi.org/10.55562/jrucs.v59i1.19>

### 1. Introduction

Health crisis prediction - predicting high-risk clinical trajectories just hours before their overt decompensation - is one of the most successful interventions in our current intensive care system [1]. In the Iraqi context, the ratio of doctors to patients in Baghdad's public teaching

hospitals is about 1:450 during peak demand times, and because of this, the delay in diagnosis of deteriorating patients has serious consequences [2]. The Iraqi Ministry of Health (2023) reported that 34.7% of preventable ICU admissions due to delayed diagnosis of sepsis, acute heart failure exacerbations, and pulmonary deterioration are reported in clinical records up to 6-12 hours ahead of time [3]. Despite the rapid development of machine learning and data mining techniques, they still face many challenges in applying them in real-life Iraqi medical systems. The first is the need for centralized data mining. This involves the aggregation of protected health information (PHI) from different institutions into a single computational environment, which is incompatible with Iraqi Law No. 65 of 2015, which requires personal data protection and is based on the ethical requirements of the Iraqi Medical Syndicate. The second challenge is complying with the General Data Protection Regulation (GDPR) and research funding bodies [4]. The second issue that is usually not mentioned in the literature is that centralized models trained on pooled multi-hospital data tend to be biased towards different demographic groups. For example, in three major hospital datasets from Baghdad, standard gradient boosting models have an F1 score on average that is 19.3 percentage points lower for patients over 70 and 14.7 percentage points lower for female patients than male patients from 30 to 55, which are the largest demographic groups in the data [5]. Federated learning (FL), introduced by McMahan et al. [6], is one of the most significant changes in model training. The model is trained on multiple nodes, and only the parameters of the model are updated from a central server, rather than just sharing the raw data. However, traditional FL algorithms, such as FedAvg and its derivatives, were primarily designed for applications in computer vision and natural language processing, making them less suitable for the heterogeneous, longitudinal, and clinically coded data structures found in hospital electronic health records (EHRs). Additionally, standard FL frameworks do not offer formal guarantees of demographic fairness. Recent studies indicate that FedAvg can exacerbate existing dataset-level biases when the data distribution is not independent and identically distributed (i.i.d.) across different hospital silos [7].

We are providing the following original and verifiable solutions to address the gaps described in this paper:

- 1) FedFairMine Architecture: A complete federated data mining system tailored to structured clinical EHR data, incorporating asynchronous silo communication adapted to the intermittent connectivity constraints of Baghdad hospital networks.
- 2) Fairness-Constrained Aggregation: A novel dual-objective optimization algorithm, FedFair-Agg, that enforces Equalized Odds constraints across demographic groups during global model aggregation, without requiring access to raw demographic labels at the central server.
- 3) Calibrated Differential Privacy: A  $(\epsilon, \delta)$ -differentially private mechanism with  $\epsilon = 0.8$  that considers the sensitivity of Iraqi clinical EHR feature vectors to balance privacy budget with predictive utility.
- 4) Empirical Baghdad validation: Based on 47,312 real patient episodes in six teaching hospitals in Baghdad, our study is the first research that has been carried out in the field of federated health crisis prediction based on Iraqi clinical data. In Section II, related works are reviewed. In Section III, the proposed FedFairMine methodology is presented, and mathematical ideas are discussed. In section IV, the experimental setup is presented in Baghdad hospitals. Section V contains quantitative results and comparisons. In section VI, we discuss the social impact and the ethical position of people, and the relevance of the SDGs. The limitations of Section VII are summarized here.

### 1.1. Data Mining for Healthcare Crisis Prediction

Data mining in clinical early-warning systems has been extensively used in the past decades and still faces many challenges. Ren et al. [8] designed DeepSOFA: a recurrent neural network that uses time series data from the ICU of MIMIC-III and achieves an area under the receiver operating characteristic curve (AUROC) of 0.869 in sepsis detection. However, when applied to non-Western patients, DeepSOFA suffers from a severe performance drop and only produces an AUROC of 0.12. This is particularly seen in Iraq, where comorbidity profiles are not only more common than

those of North American patients but also include a higher prevalence of type 2 diabetes and chronic kidney disease [9].

The gradient boosting decision tree (GBDT) is still the most popular method for structured EHR prediction tasks. Chen and colleagues [10] compared XGBoost with logistic regression and random forests for 23 clinical classification tasks. However, the models are trained using data from three medical centers in the United States, and there is concern about generalization in the other areas of the country. Moor et al. [11] proposed an enhancement of the Early Warning Score using Transformer architectures on eICU and PhysioNet 2019 datasets. But their approach requires GPU clusters for both training and inference resources that are largely unavailable in Baghdad hospitals, where CPU-only servers are the standard computing platforms. Most critically, there has been no study that has developed a validated health crisis prediction pipeline for Iraqi hospitals. With a population of 43 million and more than 35 teaching hospitals in Baghdad alone, this represents a significant geographic and demographic gap in the global health AI literature.

### **1.2. Federated Learning in Healthcare**

The application of federated learning (FL) to clinical prediction has seen significant growth since 2021. Rieke et al. [12] introduced a comprehensive framework for privacy-preserving collaborative learning in medical imaging. Their research demonstrated that federated models trained across 20 international sites can achieve centralized performance in brain tumor segmentation, with a DICE score deviation of only 3%. Similarly, Tahir et al. [13] applied this methodology to predict COVID-19 outcomes across 20 global institutions, finding that federated models were able to match or even surpass the performance of locally trained models at 14 out of the 20 sites. However, it is important to note that both studies focused on imaging data, where pixel-level privacy risks differ from those associated with structured electronic health record (EHR) data mining. In EHR data, individual fields such as diagnosis codes and medication lists are more readily identifiable.

Pfohl et al. [14] demonstrated that federated models trained using FedAvg on electronic health record (EHR) data exhibit significant performance disparities across client sites, with F1-score variances reaching up to 0.23. Such variability is unacceptable in safety-critical healthcare applications. To address the heterogeneity of EHR data in silos, Nguyen et al. [15] present FedHealth with personalized layers. But they lack mechanisms for fairness enforcement and no privacy protection yet beyond secure aggregation. This highlights a significant methodological limitation: no federated learning (FL) framework currently addresses (a) the heterogeneity of EHR data across hospital silos; (b) formal differential privacy guarantees for clinical data sensitive to the data; and (c) demographic fairness enforcement in non-i.i.d. data cases. FedFairMine is intended to fill this gap.

### **1.3. Algorithmic Fairness in Federated Systems**

Algorithmic fairness has emerged as an important area of research on federated learning, with the findings of Caton and Haas. [16], who found systematic bias amplification in distributed training. In response to this challenge, Li et al. [17] proposed q-FedAvg, which simplifies the loss function for underperforming clients and thus improves performance, although it does not address demographic parity. Papadaki et al. [18] developed a minimax fairness formulation for federated learning to reduce the worst-group error, most effectively in tabular data. However, they failed to account for demographic and geographic heterogeneity in hospitals' populations, a crucial issue in multi-silo hospital networks, as geographic differences always obscure demographic differences.

Du et al. [19] proposed a fairness-aware federated optimization approach that considers agnostic training goals. Ezzeldin et al. [20] considered fairness in healthcare-focused federated learning (FL) and demonstrated that standard FL algorithms can generate models that are generally accurate for minority patient subgroups, but fair in general. Their investigation into racial and gender disparities in federated sepsis prediction found that inequality in Equalized Odds is still

present after post-hoc debiasing of these algorithms. Such a finding underscores the necessity of developing an in-training fairness constraint.

#### 1.4. Differential Privacy in Federated Learning

The study of privacy in federated systems started with Abadi et al. [21], who introduced DP-SGD to guarantee  $(\epsilon, \delta)$ -differential privacy via gradient clipping and Gaussian noise injection. Building on this foundation, McMahan et al. [22] extended these ideas to federated systems by providing user-level differential privacy, which can be applied to millions of clients. More recently, Agarwal et al. [23] proposed adaptive clipping mechanisms that dynamically adjust noise calibration based on the distributions of gradient norms. This approach significantly reduces the accuracy loss due to clipping, with a 6.2% improvement in accuracy compared to fixed threshold clipping with the same privacy budget. One serious limitation of DP-FL (Differentially Private Federated Learning) approaches is that they assume that all client datasets come from the same distribution. This is a major problem in hospital networks where patient characteristics can vary significantly by geographic catchment area. Girgis et al. [24] find a shuffling solution to this problem, but their work does not consider fairness. Our proposed calibrated DP mechanism complements this work by effectively balancing privacy expense against demographic fairness.

#### 1.5. Research Gap Synthesis

In the last section, we present a large and largely unexplored area of research that is: (i) federated data mining architecture that can handle the heterogeneity of structured Electronic Health Records (EHRs), (ii) formally verified differential privacy for the sensitive nature of clinical data, (iii) enforcement of demographic fairness during training with Equalized Odds constraints and (iv) empirical validation in an urban hospital network in a developing nation. Table I compares FedFairMine to 12 previous studies on these dimensions.

**Table I: Comparative Analysis of Related Works — Coverage of Key Dimensions**

Work / Framework	FL Arch.	EHR Data	Diff. Privacy	Fairness	Dev. Nation	Health Crisis
Rieke et al. [12]	✓	✗	Partial	✗	✗	✗
Dayan et al. [13]	✓	✗	✗	✗	✗	✓
Pfohl et al. [14]	✓	✓	✗	Partial	✗	✗
Nguyen et al. [15]	✓	✓	✗	✗	✗	✓
Li et al. [17]	✓	✗	✗	Partial	✗	✗
Papadaki et al. [18]	✓	Partial	✗	✓	✗	✗
Ezzeldin et al. [20]	✓	✓	✗	✓	✗	✓
Agarwal et al. [23]	✓	✗	✓	✗	✗	✗
<b>FedFairMine (Ours)</b>	✓	✓	✓	✓	✓	✓

## 2. Proposed Methodology — Fedfairmine Framework

### 2.1. System Architecture Overview

FedFairMine is organized into three layers: (1) Hospital Edge Nodes (HENs), located in a teaching hospital in Baghdad, (2) Secure Aggregation Server (SAS) at the Iraqi Ministry of Health data center, (3) Global Model Repository that provides secure API access to all participating HENs. The central point of this architecture is that no patient data, such as features, diagnoses, medications, demographics, etc., leaves the local hospital network. In Figure 1, we show the end-to-end architecture of our system. Each Health Engagement Node (HEN) is provided with its own local training pipeline of (a) the EHR preprocessing engine, (b) the local feature extractor (LFE), (c) a task-specific predictive head, and (d) a differential privacy noise injection module. The Secure Aggregation Server (SAS) receives only encrypted, noise-injected gradient updates. It performs weighted aggregation while adhering to fairness constraints and subsequently transmits the updated global model back to all HENs for a total of  $T = 120$  global communication rounds. Finally, the global model is delivered to each hospital's clinical decision support interface.

## 2.2. Local Feature Extraction Engine

EHR data in Baghdad hospitals is stored in a variety of formats, from structured SQL tables (ICD-10 diagnoses, laboratory results) to unstructured physician notes in Arabic and semi-structured nursing observation charts. The Local Feature Extraction Engine (LFEE) generates this together into a unique clinical feature vector  $\varphi_i \in \mathbb{R}^{128}$  for each patient episode  $i$ . The feature vector comprises six clinically validated sub-domains: (1) Vital Signs Trajectory (32 dimensions): mean, standard deviation, minimum, maximum and rate-of-change over 1h, 6h and 12h windows for heart rate, systolic/diastolic blood pressure, SpO<sub>2</sub>, respiratory rate and temperature; (2) Laboratory Panel (28 dimensions): full blood count, metabolic panel, coagulation studies, lactate, procalcitonin, and BNP; (3) Medication Exposure (16 dimensions): binary indicators for 16 high-risk medication classes (vasopressors, anticoagulants, diuretics etc.); (4) Comorbidity Burden (18 dimensions): CCI components are written as binary flags; (5) Temporal Admission Context (16 dimensions): time-of-admission, day-of-week, season, ward type; (6) Socioeconomic Indicators (18 dimensions): district-level deprivation index based on Iraqi census data, insurance status, referral pathway. Most importantly, the LFEE computes each feature locally. Only the resulting 128-dimensional vector  $\varphi_i$  (after privacy perturbation) is used for training the local model. The Python implementation of LFEE is described in Algorithm 1 (Code Block 1).

### Code Block 1: Local Feature Extraction Engine (LFEE) — Python Implementation

```
# FedFairMine — Local Feature Extraction Engine (LFEE)
# Runs entirely on each Hospital Edge Node (HEN) — no raw data leaves the hospital

import numpy as np
import pandas as pd
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.impute import KNNImputer
from scipy.stats import zscore
import warnings
warnings.filterwarnings('ignore')

class LocalFeatureExtractor:
    """
    Privacy-safe local feature extraction for EHR crisis prediction.
    All computation happens on the local Hospital Edge Node (HEN).
    Only the resulting feature matrix — never raw EHR rows — is used
    for downstream federated training.
    """

    VITAL_COLS = ['hr', 'sbp', 'dbp', 'spo2', 'rr', 'temp']
    WINDOWS = [1, 6, 12] # hours
    LAB_COLS = ['wbc', 'hgb', 'plt', 'na', 'k', 'cr', 'bun',
                'glucose', 'lactate', 'pct', 'bnp', 'inr',
                'alt', 'ast', 'bilirubin', 'albumin',
                'ph', 'pco2', 'po2', 'hco3',
                'troponin', 'CRP', 'ferritin',
                'fibrinogen', 'd-dimer', 'procalcitonin',
                'ldh', 'il6']
    MED_CLASSES = ['vasopressor', 'anticoag', 'diuretic', 'abx',
                  'steroid', 'antifungal', 'antiviral', 'insulin',
                  'sedation', 'analgesia', 'proton_pump',
                  'bronchodilator', 'antiarrhythmic',
                  'acei_arb', 'statin', 'heparin']

    def __init__(self, scaler_path=None):
        self.imputer = KNNImputer(n_neighbors=5, weights='distance')
        self.scaler = StandardScaler()
```

```

self._fitted = False
if scaler_path:
    self._load_scaler(scaler_path)

# — 1. Vital-sign trajectories (32 features) —————
def _extract_vitals (self, df_vitals: pd.DataFrame) -> np.ndarray:
    features = []
    for col in self.VITAL_COLS:
        col_data = df_vitals[col].dropna().values
        if len(col_data) == 0:
            features.extend([0.0]*5)
            continue
        features.append(float(np.mean(col_data)))
        features.append(float(np.std(col_data) if len(col_data)>1 else 0))
        features.append(float(np.min(col_data)))
        features.append(float(np.max(col_data)))
        # Rate-of-change: last value minus first / elapsed time
        roc = (col_data[-1] - col_data[0]) / max(len(col_data)-1, 1)
        features.append(float(roc))
    return np.array(features[:32], dtype=np.float32)

# — 2. Laboratory panel (28 features) —————
def _extract_labs (self, df_labs: pd.DataFrame) -> np.ndarray:
    vec = np.zeros(28, dtype=np.float32)
    for idx, lab in enumerate(self.LAB_COLS):
        if lab in df_labs.columns:
            val = df_labs[lab].dropna().values
            vec[idx] = float(np.mean(val)) if len(val) > 0 else 0.0
    return vec

# — 3. Medication exposure (16 binary features) —————
def _extract_medications (self, med_list: list) -> np.ndarray:
    vec = np.zeros(16, dtype=np.float32)
    for idx, cls in enumerate(self.MED_CLASSES):
        vec[idx] = 1.0 if any(cls in m.lower() for m in med_list) else 0.0
    return vec

# — 4. Comorbidity Charlson Index (18 features) —————
def _extract_comorbidities (self, icd_codes: list) -> np.ndarray:
    CCI_MAP = {
        'E11': 0, 'I50': 1, 'I21': 2, 'J44': 3, 'N18': 4,
        'K70': 5, 'C': 6, 'I10': 7, 'J18': 8, 'A41': 9,
        'G30': 10, 'F20': 11, 'M05': 12, 'K57': 13, 'G35': 14,
        'I63': 15, 'K92': 16, 'B18': 17
    }
    vec = np.zeros(18, dtype=np.float32)
    for code in icd_codes:
        for prefix, idx in CCI_MAP.items():
            if code.startswith(prefix):
                vec[idx] = 1.0
    return vec

# — 5. Temporal context (16 features) —————
def _extract_temporal(self, admission_dt) -> np.ndarray:
    import math
    h = admission_dt.hour

```

```

d = admission_dt.dayofweek
vec = np.array([
    math.sin(2*math.pi*h/24), math.cos(2*math.pi*h/24),
    math.sin(2*math.pi*d/7), math.cos(2*math.pi*d/7),
    float(admission_dt.month in [12,1,2]), # winter
    float(admission_dt.month in [3,4,5]), # spring
    float(admission_dt.month in [6,7,8]), # summer
    float(admission_dt.month in [9,10,11]), # autumn
    float(h >= 8 and h < 20), # daytime
    float(d < 5), # weekday
    0.0, 0.0, 0.0, 0.0, 0.0, 0.0 # ward-type placeholders
], dtype=np.float32)
return vec

```

# — 6. Socioeconomic indicators (18 features) —————

```

def _extract_socioeconomic(self, district_id: int,
                          insurance: int, referral: str) -> np.ndarray:
    # District deprivation index from Iraqi census 2023 (0-1 scale)
    DISTRICT_DEP = {1:0.72, 2:0.65, 3:0.81, 4:0.55, 5:0.69,
                    6:0.78, 7:0.62, 8:0.74, 9:0.58, 10:0.83}
    dep = DISTRICT_DEP.get(district_id, 0.70)
    ref_enc = {'emergency':1, 'gp':0, 'specialist':0, 'self':0}
    vec = np.zeros(18, dtype=np.float32)
    vec[0] = dep
    vec[1] = float(insurance)
    vec[2] = ref_enc.get(referral.lower(), 0)
    return vec

```

# — Master extraction pipeline —————

```

def extract(self, patient_episode: dict) -> np.ndarray:
    vitals = self._extract_vitals(patient_episode['vitals_df'])
    labs = self._extract_labs(patient_episode['labs_df'])
    meds = self._extract_medications(patient_episode['medications'])
    comor = self._extract_comorbidities(patient_episode['icd_codes'])
    temporal = self._extract_temporal(patient_episode['admission_dt'])
    socio = self._extract_socioeconomic(
        patient_episode['district_id'],
        patient_episode['insurance'],
        patient_episode['referral'])
    phi = np.concatenate([vitals, labs, meds, comor, temporal, socio])
    assert phi.shape[0] == 128, f'Expected 128 features, got {phi.shape[0]}'
    return phi.astype(np.float32)

```

```

def fit_transform(self, episodes: list) -> np.ndarray:
    X = np.stack([self.extract(ep) for ep in episodes])
    X = self.imputer.fit_transform(X)
    X = self.scaler.fit_transform(X)
    self._fitted = True
    return X

```

```

def transform(self, episodes: list) -> np.ndarray:
    assert self._fitted, 'Call fit_transform first on local data'
    X = np.stack([self.extract(ep) for ep in episodes])
    X = self.imputer.transform(X)
    X = self.scaler.transform(X)
    return X

```

### 2.3. Privacy-Preserving Mechanism: Calibrated Differential Privacy

We implement  $(\epsilon, \delta)$ -differential privacy via the Gaussian mechanism applied to per-sample gradient contributions during local model training. Formally, the privacy guarantee is defined as: for any two adjacent datasets  $D$  and  $D'$  differing in exactly one patient record, and for any measurable output set  $S$ , the mechanism  $M$  satisfies:

$$\Pr[M(D) \in S] \leq e^{\epsilon} \cdot \Pr[M(D') \in S] + \delta \quad (1)$$

We set  $\epsilon = 0.8$  and  $\delta = 1 \times 10^{-5}$ , calibrated to ensure that no individual patient record has a disproportionate influence on the communicated gradient update. The noise scale  $\sigma$  is computed via:

$$\sigma = (C \cdot \sqrt{2 \cdot \ln(1.25/\delta)}) / \epsilon \quad (2)$$

where  $C$  is the L2 norm clipping threshold, set to  $C = 1.2$  based on empirical analysis of gradient norm distributions across our six Baghdad hospital datasets. The implementation in Code Block 2 details the complete DP-SGD training loop.

#### Code Block 2: Differential Privacy Training Loop (DP-SGD) with Gradient Clipping

```
# FedFairMine — Differential Privacy Module
# ( $\epsilon=0.8$ ,  $\delta=1e-5$ ) — Gaussian Mechanism with per-sample gradient clipping

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader, TensorDataset
import numpy as np
from typing import Tuple

class ClinicalCrisisPredictor(nn.Module):
    """Lightweight MLP optimized for CPU inference on Baghdad hospital servers."""
    def __init__(self, input_dim=128, hidden=[256,128,64], dropout=0.3):
        super().__init__()
        layers = []
        in_d = input_dim
        for h in hidden:
            layers += [nn.Linear(in_d, h), nn.BatchNorm1d(h),
                       nn.ReLU(), nn.Dropout(dropout)]
            in_d = h
        layers.append(nn.Linear(in_d, 1))
        self.net = nn.Sequential(*layers)

    def forward(self, x):
        return self.net(x).squeeze(-1)

class DPTrainer:
    """
    DP-SGD trainer with per-sample gradient clipping and Gaussian noise.
    Provides ( $\epsilon=0.8$ ,  $\delta=1e-5$ )-DP guarantee per communication round.
    """
    def __init__(self, model: nn.Module, epsilon=0.8, delta=1e-5,
                 clip_norm=1.2, local_epochs=5, lr=1e-3, batch_size=64):
        self.model = model
        self.epsilon = epsilon
        self.delta = delta
        self.C = clip_norm
        self.local_epochs = local_epochs
```

```

self.lr      = lr
self.batch_size = batch_size
# Gaussian mechanism noise scale
self.sigma = (clip_norm * np.sqrt(2 * np.log(1.25 / delta))) / epsilon
self.optimizer = optim.Adam(model.parameters(), lr=lr, weight_decay=1e-4)
self.criterion = nn.BCEWithLogitsLoss(reduction='none')

def _per_sample_gradients(self, X_batch, y_batch):
    """Compute per-sample gradients for clipping — core DP-SGD step."""
    per_sample_grads = []
    for xi, yi in zip(X_batch, y_batch):
        self.optimizer.zero_grad()
        out = self.model(xi.unsqueeze(0))
        loss = self.criterion(out, yi.unsqueeze(0).float())
        loss.backward()
        grads = [p.grad.clone() for p in self.model.parameters()
                 if p.grad is not None]
        per_sample_grads.append(grads)
    return per_sample_grads

def _clip_and_aggregate(self, per_sample_grads):
    """Clip per-sample gradients by L2 norm, then sum."""
    agg = None
    for grads in per_sample_grads:
        total_norm = torch.sqrt(sum(g.norm()**2 for g in grads))
        clip_factor = min(1.0, self.C / (total_norm.item() + 1e-8))
        clipped = [g * clip_factor for g in grads]
        if agg is None:
            agg = clipped
        else:
            agg = [a + c for a, c in zip(agg, clipped)]
    return agg

def _add_gaussian_noise(self, aggregated_grads, n_samples):
    """Add calibrated Gaussian noise to satisfy  $(\epsilon, \delta)$ -DP."""
    return [
        g + torch.randn_like(g) * self.sigma * self.C
        for g in aggregated_grads
    ]

def local_train(self, X: np.ndarray, y: np.ndarray) -> dict:
    """
    Run DP-SGD local training. Returns state_dict for aggregation.
    X: (N, 128) feature matrix | y: (N,) binary crisis labels
    """
    X_t = torch.tensor(X, dtype=torch.float32)
    y_t = torch.tensor(y, dtype=torch.float32)
    dataset = TensorDataset(X_t, y_t)
    loader = DataLoader(dataset, batch_size=self.batch_size, shuffle=True)

    self.model.train()
    for epoch in range(self.local_epochs):
        for X_b, y_b in loader:
            per_grads = self._per_sample_gradients(X_b, y_b)
            agg_grads = self._clip_and_aggregate(per_grads)
            noisy_grads = self._add_gaussian_noise(agg_grads, len(X_b))

```

```

# Apply noisy aggregated gradient update
self.optimizer.zero_grad()
for param, grad in zip(
    [p for p in self.model.parameters() if p.requires_grad],
    noisy_grads):
    param.grad = grad / len(X_b)
self.optimizer.step()

return self.model.state_dict()

@property
def privacy_budget_used(self) -> Tuple[float, float]:
    return (self.epsilon, self.delta)

```

#### 2.4. Fairness-Constrained Federated Aggregation — FedFair-Agg

The core innovation of FedFairMine is the FedFair-Agg aggregation algorithm, which extends standard weighted averaging with an explicit fairness regularization term. Let  $K$  denote the number of hospital silos,  $n_k$  the local dataset size of silo  $k$ , and  $G = \{g_1, g_2, \dots, g_M\}$  a set of  $M$  demographic groups (e.g., age cohorts, gender, socioeconomic quartiles). The global model parameters  $\theta^*$  are obtained by solving:

$$\theta = \underset{\theta}{\operatorname{argmin}} \left\{ \sum_{k=1}^K (n_k / N) \mathcal{L}_k(\theta) + \lambda \sum_{m=1}^M \sum_{l \neq m} |EOD(g_m, g_l; \theta)|^2 \right\} \quad (3)$$

Where  $\mathcal{L}_k(\theta)$  is the binary cross-entropy loss on silo  $k$ 's local validation set,  $N = \sum n_k$  is the total number of patients,  $EOD(g_m, g_l; \theta)$  is the Equalized Odds Difference between groups  $g_m$  and  $g_l$  (the maximum difference of TPR and FPR value across groups) and  $\lambda = 0.15$  is the fairness regularization weight due to cross-validation. The SAS does not access the data for demographic labels. Each HEN only sends its local EOD estimates, which are calculated locally on demographic-stratified validation splits, as well as the encrypted gradient update. The SAS weights the aggregation inversely to the number of reported EOD violations and penalizes model updates with a higher demographic disparity. Code Block 3 contains the entire FedFair-Agg implementation.

#### Code Block 3: FedFair-Agg — Fairness-Constrained Federated Aggregation Algorithm

```

# FedFairMine — FedFair-Agg: Fairness-Constrained Aggregation
# Runs on the Secure Aggregation Server (SAS) — Ministry of Health Data Center

import numpy as np
import torch
from collections import OrderedDict
from typing import List, Dict, Tuple
import copy

def equalized_odds_difference(y_true: np.ndarray, y_pred: np.ndarray,
                             sensitive: np.ndarray) -> float:
    """
    Compute Equalized Odds Difference (EOD) between two demographic groups.
    EOD = max(|TPR_0 - TPR_1|, |FPR_0 - FPR_1|)
    Groups are binary-encoded (0/1) in `sensitive`.
    """
    groups = np.unique(sensitive)
    if len(groups) < 2:
        return 0.0

    def group_rates(g):
        mask = sensitive == g
        yt, yp = y_true[mask], y_pred[mask]

```

```

if yt.sum() == 0 or (1-yt).sum() == 0:
    return 0.0, 0.0
tpr = ((yp == 1) & (yt == 1)).sum() / (yt == 1).sum()
fpr = ((yp == 1) & (yt == 0)).sum() / (yt == 0).sum()
return float(tpr), float(fpr)

```

```

tpr0, fpr0 = group_rates(groups[0])
tpr1, fpr1 = group_rates(groups[1])
return max(abs(tpr0 - tpr1), abs(fpr0 - fpr1))

```

```
class FedFairAgg:
```

```
    """
```

```

    Privacy-preserving, fairness-constrained federated aggregation.
    The server receives only (state_dict, n_samples, eod_reports) — no raw data.
    """

```

```

def __init__(self, global_model: torch.nn.Module,
             lambda_fair: float = 0.15,
             n_rounds: int = 120):
    self.global_model = copy.deepcopy(global_model)
    self.lambda_fair = lambda_fair
    self.n_rounds = n_rounds
    self.round_history = []

```

```

def _compute_fairness_penalty(self,
                             eod_reports: List[Dict]) -> np.ndarray:
    """
    Each HEN reports EOD per demographic axis:
    eod_reports = [{age: 0.03, gender: 0.05, socio: 0.04}, ...]
    Returns penalty weight per client (higher EOD → lower weight).
    """
    eod_scalars = np.array([
        np.mean(list(report.values())) for report in eod_reports
    ])
    # Inverse penalty: clients with lower EOD get higher weight
    fair_weights = 1.0 / (1.0 + self.lambda_fair * eod_scalars)
    return fair_weights / fair_weights.sum()

```

```

def aggregate(self,
             client_updates: List[Tuple[OrderedDict, int, Dict]],
             ) -> OrderedDict:
    """
    Aggregate client model updates with fairness constraints.

    Args:
        client_updates: list of (state_dict, n_samples, eod_report) tuples
    Returns:
        Aggregated global model state_dict
    """
    state_dicts = [u[0] for u in client_updates]
    n_samples = np.array([u[1] for u in client_updates], dtype=float)
    eod_reports = [u[2] for u in client_updates]

    # Standard FedAvg weights (proportional to dataset size)
    fedavg_weights = n_samples / n_samples.sum()

```

```

# Fairness penalty weights (inverse to EOD magnitude)
fair_weights = self._compute_fairness_penalty(eod_reports)

# Combine: joint weight balances data size and fairness penalty
joint_weights = 0.7 * fedavg_weights + 0.3 * fair_weights
joint_weights /= joint_weights.sum() # renormalize

# Weighted parameter aggregation
global_state = OrderedDict()
for key in state_dicts[0].keys():
    stacked = torch.stack([
        sd[key].float() * w
        for sd, w in zip(state_dicts, joint_weights)
    ])
    global_state[key] = stacked.sum(dim=0)

self.global_model.load_state_dict(global_state)

# Log round statistics
mean_eod = np.mean([
    np.mean(list(r.values())) for r in eod_reports
])
self.round_history.append({
    'mean_eod' : mean_eod,
    'fair_weights': fair_weights.tolist(),
    'joint_weights': joint_weights.tolist()
})

return global_state

def run_federation(self,
    hospital_trainers: list,
    hospital_datasets: list) -> torch.nn.Module:
    """
    Main federated training loop.
    hospital_trainers: list of DPTrainer instances (one per HEN)
    hospital_datasets: list of (X_train, y_train, demo_df) per HEN
    """
    for round_t in range(self.n_rounds):
        client_updates = []
        # Distribute global model to all HENs
        global_sd = copy.deepcopy(self.global_model.state_dict())

        for trainer, (X, y, demo) in zip(hospital_trainers, hospital_datasets):
            trainer.model.load_state_dict(global_sd)
            # Local DP training — no raw data leaves hospital
            local_sd = trainer.local_train(X, y)

            # Compute local EOD reports per demographic axis
            with torch.no_grad():
                logits = trainer.model(torch.tensor(X, dtype=torch.float32))
                y_pred = (torch.sigmoid(logits) > 0.5).numpy().astype(int)

            eod_report = {
                for axis in ['age_group', 'gender', 'socio_quartile']:
                    if axis in demo.columns:

```

```

unique_g = demo[axis].unique()
if len(unique_g) >= 2:
    sens = (demo[axis] == unique_g[0]).values.astype(int)
    eod_report[axis] = equalized_odds_difference(
        y, y_pred, sens)

client_updates.append((local_sd, len(X), eod_report))

# Server aggregates with fairness constraints
self.aggregate(client_updates)

if (round_t + 1) % 20 == 0:
    mean_eod = self.round_history[-1]['mean_eod']
    print(f'Round {round_t+1:3d}/{self.n_rounds} | ',
          f'Mean EOD: {mean_eod:.4f}')

return self.global_model

```

### 3. Experimental Setup - Baghdad Hospitals Case Study

#### 3.1. Participating Institutions and Ethical Approval

The ethics of this study were approved by the Iraqi National Health Research Committee (Approval Ref: NHRC-2023-0417) and the Institutional Review Boards of the six participating hospitals. Data collection was conducted in accordance with Iraqi Law No. 65/2015 on personal data protection. The patient records were de-identified at the source using the Safe Harbor method before any computational processing, and their analysis was independently conducted by the hospital data protection officer at each site. The six hospitals involved are a very diverse group of hospitals in Baghdad and are geographically diverse: Al-Kindy Teaching Hospital (Rusafa, 720 beds); Ghazi Al-Hariri Hospital for Surgical Specialties (Medical City Complex, 400 beds); Ibn Sina Teaching Hospital (Mansour, 650 beds); Baghdad Medical City (Central Baghdad, 1,200 beds); Al-Yarmouk Teaching Hospital (Karkh, 580 beds); Al-Kadhimiya Teaching Hospital (Al-Kadhimiya, 510 beds). These institutions collectively provide about 2.8 million annual patient visits and are high-acuity tertiary referral centres and mixed-acuity district hospitals.

#### 3.2. Dataset Characteristics and Crisis Definition

The study data contains 47,312 de-identified patient episodes from January 2021 to December 2024. Each episode is a complete hospital admission, including all recorded vitals, lab results, medication administration, physician orders, and discharge diagnosis. The main prediction goal is the development of a health crisis event (label = 1), which can be (a) unplanned ICU transfer within 24 hours, (b) cardiac arrest in the hospital, (c) rapid response team activation, or (d) unexpected mechanical ventilation initiation. This definition was confirmed by five senior intensivists from the study sites. Of the 47,312 patient episodes reviewed in the study, 8,234 (17.4%) were identified as a health crisis. Table II shows the distribution of data between hospitals and demographic groups, highlighting the careful demographic heterogeneity that is necessary to test fairness mechanisms.

**Table II: Dataset Statistics Across Six Baghdad Hospital Nodes**

Hospital	Episodes	Crisis %	Age >65 %	Female %	Low SES %	Missing Data %
Al-Kindy TH	8,412	19.2%	28.4%	42.1%	61.3%	14.7%
Ghazi Al-Hariri	6,891	15.8%	22.7%	38.9%	44.2%	11.2%
Ibn Sina TH	8,103	18.1%	31.2%	45.6%	52.8%	16.4%
Baghdad Medical City	12,204	16.4%	26.8%	41.3%	38.7%	9.8%
Al-Yarmouk TH	6,918	17.9%	29.1%	43.8%	57.4%	13.1%
Al-Kadhimiya TH	4,784	18.7%	33.4%	47.2%	63.1%	15.8%
<b>Total / Mean</b>	<b>47,312</b>	<b>17.4%</b>	<b>28.6%</b>	<b>43.2%</b>	<b>52.9%</b>	<b>13.5%</b>

### 3.3. Data Preprocessing Pipeline

Real-world EHR data in Iraqi hospitals have some of the quality issues: laboratory values missing (mean 13.5%, range 9.8%–16.4% across sites), duplicate entries from multi-ward transfers, timing issues between nursing and physician records, and coding discrepancies between different ICD-10-CM and ICD-10-AM systems. Our preprocessing pipeline addresses these issues by implementing a standard five-stage process that is executed in each hospital’s own environment.

**Deduplication:** The episode with the same admission/discharge time and patient pseudo-IDs is merged, retaining the observation with the most complete laboratory panel. **Stage 2: Temporal matching:** The data is re-sampled to 1h bins for each hospital in terms of its admission time by forward-fill, and this is done with a time lag of up to 4h. **Stage 3: Missing value Imputation:** We apply KNN imputation ( $k=5$ ) independently within each hospital's dataset. No inter-hospital imputation is performed, keeping the data relatively local. **Stage 4: Outlier Analysis:** The most unlikely values ( $SpO_2 > 100\%$ , heart rate  $> 350$  bpm) are discarded, and NaN is substituted for the ones that are too high to be sent to the lab for the imputation, and the data points are marked with hospital-specific physiological boundary tables by clinical staff. **Stage 5: Correct for class imbalance.** As the 17.4% crisis prevalence has been reported in this study, SMOTE-Tomek oversampling is used locally in each hospital training area to generate a subset of minority-class episodes that are not marked as borderline. Code Block 4 presents the complete preprocessing pipeline and the SMOTE-Tomek implementation adapted to structured EHR data.

#### Code Block 4: EHR Preprocessing Pipeline with SMOTE-Tomek Resampling

```
# FedFairMine — Local EHR Preprocessing Pipeline
# Executed entirely within each hospital boundary (HEN)

import pandas as pd
import numpy as np
from sklearn.impute import KNNImputer
from imblearn.combine import SMOTETomek
from imblearn.over_sampling import SMOTE
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import warnings
warnings.filterwarnings('ignore')

class EHRPreprocessor:
    """
    Five-stage local preprocessing pipeline for Baghdad hospital EHR data.
    Handles missing data, outliers, temporal alignment, and class imbalance
    entirely within the hospital boundary — no raw data leaves the HEN.
    """

    PHYSIO_BOUNDS = {
        'hr' : (20, 350), 'sbp' : (40, 300), 'dbp' : (10, 200),
        'spo2': (50, 100), 'rr' : (4, 80), 'temp': (30, 43),
        'wbc' : (0, 200), 'hgb' : (2, 25), 'plt' : (1, 2000),
        'cr' : (0, 30), 'na' : (100,180), 'k' : (1.5, 9),
        'glucose': (20, 800), 'lactate': (0, 30), 'ph' : (6.8, 7.8),
    }

    def __init__(self, hospital_id: str):
        self.hospital_id = hospital_id
        self.imputer = KNNImputer(n_neighbors=5, weights='distance')
        self.scaler = StandardScaler()
        self.resampler= SMOTETomek(
            smote=SMOTE(k_neighbors=5, random_state=42),
```

```

    random_state=42
)

# Stage 1: Deduplication
def deduplicate(self, df: pd.DataFrame) -> pd.DataFrame:
    df = df.sort_values('admission_dt')
    df = df.drop_duplicates(subset=['patient_id', 'admission_dt'],
                           keep='first')
    # For the same patient, keep admission with the most complete lab panel
    lab_cols = [c for c in df.columns if c.startswith('lab_')]
    df['_completeness'] = df[lab_cols].notna().sum(axis=1)
    df = (df.sort_values('_completeness', ascending=False)
          .drop_duplicates(subset=['patient_id'], keep='first')
          .drop(columns=['_completeness']))
    return df.reset_index(drop=True)

# Stage 2: Temporal alignment (1-hour bins, 4-hour max gap forward-fill)
def align_temporal(self, df_obs: pd.DataFrame,
                  max_gap_hours: int = 4) -> pd.DataFrame:
    df_obs = df_obs.set_index('timestamp').sort_index()
    hourly = df_obs.resample('1H').mean()
    hourly = hourly.fillna(method='ffill', limit=max_gap_hours)
    return hourly.reset_index()

# Stage 3: Outlier detection and masking
def mask_outliers(self, df: pd.DataFrame) -> pd.DataFrame:
    df = df.copy()
    for col, (lo, hi) in self.PHYSIO_BOUNDS.items():
        if col in df.columns:
            out_mask = (df[col] < lo) | (df[col] > hi)
            n_out = out_mask.sum()
            if n_out > 0:
                df.loc[out_mask, col] = np.nan
                print(f' [{self.hospital_id}] {col}: {n_out} outliers masked')
    return df

# Stage 4: KNN imputation (local, within-hospital only)
def impute(self, X: np.ndarray, fit: bool = True) -> np.ndarray:
    if fit:
        return self.imputer.fit_transform(X)
    return self.imputer.transform(X)

# Stage 5: SMOTE-Tomek class balancing
def balance(self, X: np.ndarray,
           y: np.ndarray) -> tuple:
    print(f' [{self.hospital_id}] Pre-balance: '
          f'crisis={y.sum()}, non-crisis={len(y)-y.sum()}')
    X_res, y_res = self.resampler.fit_resample(X, y)
    print(f' [{self.hospital_id}] Post-balance: '
          f'crisis={y_res.sum()}, non-crisis={len(y_res)-y_res.sum()}')
    return X_res, y_res

def full_pipeline(self, df: pd.DataFrame,
                 feature_cols: list,
                 label_col: str,
                 test_size: float = 0.2
                 ) -> tuple:

```

```

"""
Run all 5 stages. Returns train/test splits, scaled and balanced.
"""
df = self.deduplicate(df)
df = self.mask_outliers(df)

X = df[feature_cols].values.astype(np.float32)
y = df[label_col].values.astype(np.int32)

X = self.impute(X, fit=True)

X_tr, X_te, y_tr, y_te = train_test_split(
    X, y, test_size=test_size, stratify=y, random_state=42)

X_tr = self.scaler.fit_transform(X_tr)
X_te = self.scaler.transform(X_te)

X_tr, y_tr = self.balance(X_tr, y_tr)

print(f'[{self.hospital_id}] Pipeline complete. '
      f'Train: {len(X_tr)}, Test: {len(X_te)}')
return X_tr, X_te, y_tr, y_te

```

### 3.4. Federated Network Configuration

The federated network was deployed across the six hospital HENs and on a star topology with the SAS running on a dedicated server in the Iraqi Ministry of Health IT infrastructure (Baghdad Al-Waziriya data center). All HEN-SAS communications were encrypted with TLS 1.3 and transmitted through a VPN tunnel designed to support this study. The HEN hardware is hospital-grade workstations (Intel Xeon E-2300 series, 32GB RAM, no GPU) representing the actual computational infrastructure in Baghdad teaching hospitals. The SAS server is an Intel Xeon Silver 4310 with 128GB RAM.  $T = 120$  global communication rounds were used for training, with every HEN performing  $E = 5$  local epochs per round in mini-batches of size  $B = 64$ . The federated system was implemented using PyTorch 2.1.0 with custom communication simulation using Python's multiprocessing library to faithfully model the latency and bandwidth constraints of Baghdad hospital network connections (average bandwidth 10 Mbps, average latency 45ms HEN-to-SAS). Table III shows the complete hyperparameter setup.

**Table III: FedFairMine Hyperparameter Configuration**

Hyperparameter	Value	Rationale
Global comm. rounds (T)	120	Convergence plateau at T=115
Local epochs (E)	5	Balance: local drift vs. convergence
Mini-batch size (B)	64	CPU memory constraint on HENs
Learning rate ( $\eta$ )	0.001	Adam with cosine decay
DP noise scale ( $\sigma$ )	Computed: $\epsilon=0.8, \delta=1e-5$	Clinical data sensitivity analysis
Gradient clip norm (C)	1.2	P95 of local gradient norms
Fairness weight ( $\lambda$ )	0.15	Cross-validated on held-out splits
Model architecture	MLP [128→256→128→64→1]	CPU-optimized depth/width
FedFair-Agg weights	0.7 FedAvg + 0.3 Fairness	Cross-validated combination
Min. participation fraction	5/6 nodes	Dropout tolerance

## 4. Results and Comparative Evaluation

### 4.1. Evaluation Metrics

The model was evaluated with a split of 80/20 train-test, which was performed independently within each hospital's local data, and all performance metrics were computed on the held-out test partitions. For the predictive outcome, we report: Accuracy, Area Under the ROC Curve (AUROC), Precision, Recall, F1-Score, and Average Precision (AP). For the fairness

assessment, we report: Equalized Odds Difference (EOD), Statistical Parity Difference (SPD), and Individual Fairness Score (IFS) computed by adding counterfactual data. The privacy expenditure is calculated as the  $(\epsilon, \delta)$  pair as a function of the RDP (Rényi Differential Privacy) accounting. Code Block 5 shows the complete evaluation framework, including the computation of fairness metrics.

### Code Block 5: Comprehensive Evaluation Framework — Performance and Fairness Metrics

```
# FedFairMine — Evaluation Framework
# Computes predictive performance AND demographic fairness metrics

import numpy as np
import pandas as pd
from sklearn.metrics import (
    accuracy_score, roc_auc_score, precision_score, recall_score,
    f1_score, average_precision_score, confusion_matrix,
    classification_report
)
from typing import Dict, Optional
import warnings
warnings.filterwarnings('ignore')

def statistical_parity_difference(y_pred: np.ndarray,
                               sensitive: np.ndarray) -> float:
    """SPD = P(Y=1|A=0) - P(Y=1|A=1)"""
    mask0 = sensitive == 0
    mask1 = sensitive == 1
    if mask0.sum() == 0 or mask1.sum() == 0:
        return 0.0
    p0 = y_pred[mask0].mean()
    p1 = y_pred[mask1].mean()
    return float(abs(p0 - p1))

def equalized_odds_difference(y_true: np.ndarray,
                              y_pred: np.ndarray,
                              sensitive: np.ndarray) -> Dict[str, float]:
    """Full EOD report: TPR gap, FPR gap, and combined EOD."""
    groups = np.unique(sensitive)
    if len(groups) < 2:
        return {'tpr_gap': 0.0, 'fpr_gap': 0.0, 'eod': 0.0}

    rates = { }
    for g in groups:
        mask = sensitive == g
        yt, yp = y_true[mask], y_pred[mask]
        tp = ((yp==1)&(yt==1)).sum()
        fp = ((yp==1)&(yt==0)).sum()
        fn = ((yp==0)&(yt==1)).sum()
        tn = ((yp==0)&(yt==0)).sum()
        tpr = tp / max(tp+fn, 1)
        fpr = fp / max(fp+tn, 1)
        rates[g] = (tpr, fpr)

    g0, g1 = groups[0], groups[1]
    tpr_gap = abs(rates[g0][0] - rates[g1][0])
    fpr_gap = abs(rates[g0][1] - rates[g1][1])
    eod = max(tpr_gap, fpr_gap)
```

```

return {'tpr_gap': tpr_gap, 'fpr_gap': fpr_gap, 'eod': eod}

def evaluate_model(y_true: np.ndarray,
                  y_prob: np.ndarray,
                  demo_df: Optional[pd.DataFrame] = None,
                  threshold: float = 0.5) -> Dict:
    """
    Comprehensive model evaluation.
    Returns a dict of performance and (optionally) fairness metrics.
    """
    y_pred = (y_prob >= threshold).astype(int)

    results = {
        'accuracy' : accuracy_score(y_true, y_pred),
        'auroc'    : roc_auc_score(y_true, y_prob),
        'precision' : precision_score(y_true, y_pred, zero_division=0),
        'recall'   : recall_score(y_true, y_pred, zero_division=0),
        'f1'       : f1_score(y_true, y_pred, zero_division=0),
        'avg_prec' : average_precision_score(y_true, y_prob),
        'threshold': threshold,
    }
    cm = confusion_matrix(y_true, y_pred)
    results['confusion_matrix'] = cm.tolist()

    if demo_df is None:
        fairness = {}
        for axis in ['age_group', 'gender', 'socio_quartile']:
            if axis in demo_df.columns:
                sensitive = (demo_df[axis] == demo_df[axis].unique()[0]).astype(int).values
                eod_res = equalized_odds_difference(y_true, y_pred, sensitive)
                spd = statistical_parity_difference(y_pred, sensitive)
                fairness[axis] = {'eod_res': eod_res, 'spd': spd}
        results['fairness'] = fairness
        # Aggregate EOD across all demographic axes
        eods = [v['eod'] for v in fairness.values()]
        results['mean_eod'] = float(np.mean(eods)) if eods else 0.0
    return results

def format_results_table(results_dict: Dict[str, Dict]) -> pd.DataFrame:
    """Format evaluation results from multiple models into a comparison table."""
    rows = []
    for model_name, res in results_dict.items():
        row = {
            'Model': model_name,
            'Accuracy': f"{res['accuracy']:.4f}",
            'AUROC' : f"{res['auroc']:.4f}",
            'Precision': f"{res['precision']:.4f}",
            'Recall' : f"{res['recall']:.4f}",
            'F1-Score': f"{res['f1']:.4f}",
            'Mean EOD': f"{res.get('mean_eod', 'N/A'):.4f}"
                       if isinstance(res.get('mean_eod'), float) else 'N/A',
        }
        rows.append(row)
    return pd.DataFrame(rows)

```

## 4.2. Comparative Performance Results

Table IV compares FedFairMine with the six baselines: (i) Centralized XGBoost trained on all data pooled (the privacy-violating upper bound); (ii) Centralized MLP; (iii) Local-Only MLP (each hospital trains separately, no federation); (iv) FedAvg [6]; (v) FedProx ( $\mu=0.01$ ) [25]; (vi) SCAFFOLD [26]. All of the federated baselines use the same ClinicalCrisisPredictor architecture as FedFairMine, but with no fairness-constrained aggregation or calibrated DP.

**Table IV: Comparative Performance and Fairness Evaluation — FedFairMine vs. Baselines**

Model	Accuracy	AUROC	Precision	Recall	F1-Score	Mean EOD↓	DP-Cert.
Centralized XGBoost†	0.923	0.961	0.887	0.901	0.894	0.087	×
Centralized MLP†	0.911	0.948	0.874	0.889	0.881	0.093	×
Local-Only MLP	0.847	0.891	0.803	0.824	0.813	0.064	×
FedAvg [6]	0.873	0.921	0.841	0.851	0.846	0.083	×
FedProx [25]	0.889	0.933	0.858	0.872	0.865	0.071	×
SCAFFOLD [26]	0.901	0.941	0.869	0.884	0.876	0.067	×
<b>FedFairMine (Ours)</b>	<b>0.918</b>	<b>0.957</b>	<b>0.903</b>	<b>0.921</b>	<b>0.912</b>	<b>0.021</b>	<b>✓</b>

† Centralized baselines violate patient privacy and are included solely as theoretical upper bounds. EOD: Equalized Odds Difference. DP-Cert.: Differential Privacy Certificate.

FedFairMine achieves an F1-score of 0.912 (far above all privacy-compliant federated baselines) and within 1.2 percentage points of the privacy-violating centralized upper bound. This small performance gap (1.2% versus the centralized MLP) shows that federated training with DP does not require a massive utility loss when properly calibrated. Most importantly, FedFairMine's Mean EOD of 0.021 is 76% lower than centralized XGBoost (0.087) and 69% lower than SCAFFOLD (0.067) as the best fairness-unaware federated baseline.

## 4.3. Fairness Analysis by Demographic Axis

Table V disaggregates the fairness metrics across the three primary demographic axes, revealing nuanced patterns in algorithmic equity. FedFairMine achieves consistent fairness improvements across all axes, with the most pronounced gain for the socioeconomic dimension (EOD reduction from 0.091 in centralized XGBoost to 0.018 in FedFairMine)—a finding of particular relevance given Baghdad's high socioeconomic stratification.

**Table V: Fairness Metrics by Demographic Axis — FedFairMine vs. Best Baseline (SCAFFOLD)**

Model	Age EOD	Age SPD	Gender EOD	Gender SPD	SES EOD	SES SPD	Mean EOD
Centralized XGBoost	0.081	0.074	0.088	0.079	0.091	0.083	0.087
SCAFFOLD	0.063	0.058	0.069	0.061	0.069	0.063	0.067
<b>FedFairMine (Ours)</b>	<b>0.019</b>	<b>0.016</b>	<b>0.025</b>	<b>0.021</b>	<b>0.018</b>	<b>0.015</b>	<b>0.021</b>

## 4.4. Convergence Behavior and Communication Efficiency

In Fig. 2, we show the convergence curves for FedFairMine and the three federated baselines over 120 communication rounds. FedFairMine converges to within 1% of its final F1-score by round 85, whereas SCAFFOLD and FedAvg converged to within 1% of each other in round 97 and 112, respectively. This is due to the fairness-weighting method, which implicitly weights the local updates from data-rich, but demographically biased silos (especially Baghdad Medical City, which contains 25.8% of the training data, but has the lowest minority-group representation), so that the global model does not overfit to majority demographic patterns early in the training process. The communication overhead per round is the 128-KB model state dict (encrypted) and 42 bytes of fairness report (3 demographic axes  $\times$  2 EOD values  $\times$  float64). Total communication cost per training run: 120 rounds  $\times$  6 clients  $\times$  128 KB = 92.2 MB, well within the capacity of Baghdad hospital network infrastructure.

#### 4.5. Privacy Expenditure Analysis

Using the Rényi Differential Privacy (RDP) accountant with  $\alpha = 10$ , the cumulative privacy expenditure in  $T = 120$  rounds with  $E = 5$  local epochs, and  $B = 64$  batch size on the largest hospital dataset ( $N = 9,763$  training samples) is  $\epsilon_{\text{RDP}} = 0.78$  at  $\delta = 1 \times 10^{-5}$ , which is within our target budget of  $\epsilon = 0.8$ . This proves that FedFairMine has completed its full training cycle within the privacy envelope, and no individual patient's data is under privacy expenditure above 0.8. That is a significant formal promise: if we observe all model updates broadcast by any HEN, we have at most  $e^{0.8} \approx 2.23$  times more information about the patient's clinical data than if we do not.

### 5. Discussion — Societal Impact and Ethical Alignment

#### 5.1. Clinical Impact on Baghdad Healthcare System

The direct clinical utility of FedFairMine is quantifiable: the model achieves a recall of 0.921 at the operating threshold (0.5) for crisis events. Of Baghdad's six participating hospitals that collectively receive 47,312 annual admissions at the complexity level studied, this means 7,589 of 8,234 crisis events are identified. We estimate that, by conservative estimates, which are based on the Iraqi Ministry of Health data, ICU admission costs  $3 \times$  the cost of a normal ward admission, and that early intervention can reduce the time in ICU by an average of 2.3 days [3]. FedFairMine will therefore reduce unnecessary or delayed ICU admissions by 22–28% for the entire system.

In addition to the economic impact, fairness guarantees also have serious equity implications. Before the FedFairMine analysis of clinical practice data of the study hospitals, the elderly patients ( $>65$ ) and those from low-SES districts in Baghdad were 23% less likely to receive timely clinical escalation than the general population, even after controlling for severity of illness, which is consistent with the prediction gaps in other models. FedFairMine's Mean EOD of 0.021 on all demographic axes means that crisis detection works well and is sensitive to the same number of patients at the same time, whether the patient is a young male worker from Al-Mansour or an elderly female resident from low-income Sadr City.

#### 5.2. Legal and Ethical Compliance

From the beginning of the study, FedFairMine was designed in accordance with the Iraqi personal data protection laws. The key compliance features are: (i) Data Minimization - only 128-dimensional feature vectors and 42-byte fairness reports are left at each hospital boundary as opposed to raw patient data; (ii) Purpose Limitation - the federated model is trained and used for health crisis prediction only in the hospital network; (iii) Formal Privacy Certification - the ( $\epsilon=0.8$ ,  $\delta=10^{-5}$ )-DP guarantee is a mathematically verifiable privacy guarantee which can be audited by the Iraqi Data Protection Authority; and (iv) Demographic Fairness Auditability. The per-axis EOD reports produced in each aggregation round are audited to ensure fairness compliance for all groups in the network. The study was fully approved by the Iraqi National Health Research Committee, and all participating hospitals are free from individual consent for receiving retrospective de-identified data, in line with international research ethics for secondary clinical data use.

#### 5.3. Alignment with UN Sustainable Development Goals

FedFairMine directly supports two UN Sustainable Development Goals that are relevant to Iraq's health development trajectory. As for SDG 3 (Good Health and Well-Being), specifically Target 3.4 (reduce non-communicable disease mortality by one third by 2030), the framework's early detection capability targets the three major causes of preventable hospital mortality in Baghdad: sepsis, acute cardiovascular decompensation, and respiratory failure. One-hour delay in sepsis treatment is estimated to increase mortality by 7% (WHO), and FedFairMine's 6-12-hour early warning window is a clinically actionable window. As for SDG 10 (Reduced Inequalities), specifically Target 10.3 (ensure equal opportunity and reduce inequality of outcome), the fairness of the framework addresses the systemic disadvantage experienced by the elderly, female, and low-socioeconomic-status patients in Baghdad's hospital system. The fact that Equalized Odds across demographics is so close to zero (0.087 to 0.021) means a real-world proof-of-concept for equitable healthcare provision in an urban developing country is achieved.

#### 5.4. Generalizability to Other Developing-Nation Hospital Networks

Although FedFairMine has been tested in Baghdad, the design of the framework can easily be generalized to the Middle East, North Africa, South Asia, and Sub-Saharan Africa with the following key characteristics: (i) heterogeneity in the EHR of hospitals; (ii) data protection laws which do not allow for centralization; (iii) a diverse population of people, which are shown to have inequitable health systems; (iv) CPU-constrained computational infrastructure. The LFEE is modular, and local feature schemas can be substituted to suit the coding system of different EHRs in the region, and FedFair-Agg is not sensitive to the demographics of the EODs.

#### 6. Conclusion, Limitations, And Future Research Directions

This study addressed the practical challenge of deploying privacy-aware, fairness-constrained federated learning in a real hospital network with limited infrastructure. The findings reported below reflect both the measurable gains of the proposed framework and the boundaries of what the current experimental setup can support.

1. FedFairMine, a privacy-preserving federated data mining framework, was developed and evaluated for early health crisis prediction across six Baghdad teaching hospitals, drawing on 47,312 patient episodes collected between 2021 and 2024.
2. The framework achieved an F1-score of 0.912 under formal differential privacy certification ( $\epsilon=0.8$ ,  $\delta=10^{-5}$ ), with an Equalized Odds Difference of 0.021 — a 76% reduction compared to the strongest centralized baseline and 69% better than the top privacy-agnostic federated baseline.
3. These results confirm that predictive accuracy, patient privacy, and demographic fairness are not mutually exclusive in federated clinical AI; the trade-offs observed in prior literature can be reduced through deliberate architectural choices.
4. The framework's performance supports its application in resource-constrained hospital networks and aligns with the equity and health-access objectives of UN SDGs 3 and 10.
5. A known computational cost remains: DP-SGD's per-sample gradient computation runs approximately  $3.8\times$  slower than standard mini-batch SGD, which is acceptable for scheduled training but may introduce delays in time-critical clinical scenarios.
6. Generalizability is limited by the dataset's geographic scope (Baghdad hospitals only) and temporal span (four years), making extrapolation to post-pandemic disease patterns or hospitals outside Iraq uncertain.
7. The absence of unstructured Arabic-language physician notes in the feature set accounts for at least part of the 1.2-point gap between FedFairMine and the centralized upper bound, a gap that structured data alone cannot.
8. Three directions are identified for follow-on work: (1) per-silo personalization layers (e.g., pFedMe or APFL) to better capture hospital-specific patient population differences without sacrificing privacy or fairness guarantees; (2) integration of an Arabic-aware clinical NLP encoder based on CAMEL-Medical to incorporate free-text physician notes, with a projected 4–7% gain in crisis prediction F1-score; and (3) asynchronous streaming federated updates that allow the global model to incorporate new clinical evidence within hours rather than days while maintaining differential privacy and fairness constraints.

#### 7. Acknowledgements

The authors are grateful to the Iraqi Ministry of Health for institutional support and data access facilitation; the medical informatics staff and nursing teams at all six participating Baghdad hospitals for clinical data curation; and the University of Baghdad College of Science's High-Performance Computing cluster for pilot experiment infrastructure. This work was not funded by external sources and was done as an independent research project under the University of Baghdad and Al-Nahrain University.

## References

- [1] Saqib M, Iftikhar M, Neha F, Karishma F, Mumtaz H. "Artificial intelligence in critical illness and its impact on patient care: a comprehensive review." *Frontiers in Medicine*. 2023; 10:1176192. <https://doi.org/10.3389/fmed.2023.1176192> [[Frontiers in Medicine]]
- [2] Ministry of Health, Republic of Iraq. Annual Statistical Report 2023: Baghdad and Governorates. Baghdad: MOH Planning Department; 2024. Available from: <https://moh.gov.iq/hci.mop.gov.iqstorage.moh.gov.iq>
- [3] R. Malaeb et al., "High mortality rates among COVID-19 intensive care patients in Iraq: insights from a retrospective cohort study at Médecins Sans Frontières supported hospital in Baghdad," *Front. Public Health*, vol. 11, p. 1185330, Aug. 2023, doi: 10.3389/fpubh.2023.1185330
- [4] M. Nahi, "Losing Citizenship by Deprivation in the Framework of Private International Law: A Comparative Study," *Journal of Lifestyle and SDGs Review*, vol. 5, no. 6, e06906, Jun. 2025. doi: 10.47172/2965-730X.SDGsReview.v5.n06.pe06906
- [5] B. Al Jorf and F. E. Shamout, "MedPatch: Confidence-Guided Multi-Stage Fusion for Multimodal Clinical Data," in *Proceedings of Machine Learning Research*, vol. 298, pp. 1-28, Aug. 2025.
- [6] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, 'Communication-efficient learning of deep networks from decentralized data,' in *Proc. 20th Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, Fort Lauderdale, FL, Apr. 2017, pp. 1273–1282.
- [7] J. Xu, B. S. Glicksberg, C. Su, P. Walker, J. Bian, and F. Wang, 'Federated learning for healthcare informatics,' *Journal of Healthcare Informatics Research*, vol. 5, no. 1, pp. 1–19, Mar. 2021.
- [8] W. Ren, J. Zhu, Z. Liu, T. Zhao, and V. Honavar, "A Comprehensive Survey of Electronic Health Record Modeling: From Deep Learning Approaches to Large Language Models," *arXiv preprint arXiv:2507.12774*, Jul. 2025. doi: 10.48550/arXiv.2507.12774
- [9] A. Asiri et al., "Clinical Profiles, Interventions, and Outcomes of Sepsis and Septic Shock in a Saudi Arabian Tertiary ICU: A Five-Year Retrospective Analysis," *Healthcare (Basel)*, vol. 14, no. 5, p. 680, Mar. 2026, doi: 10.3390/healthcare14050680.
- [10] T. Chen and C. Guestrin, 'XGBoost: A scalable tree boosting system,' in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, San Francisco, CA, Aug. 2016, pp. 785–794.
- [11] M. Moor, M. Rieck, M. Horn, C. A. Jutzeler, and K. Borgwardt, 'Early prediction of sepsis in the ICU using machine learning: A systematic review,' *Frontiers in Medicine*, vol. 8, p. 607952, Feb. 2021.
- [12] N. Rieke et al., 'The future of digital health with federated learning,' *NPJ Digital Medicine*, vol. 3, no. 1, p. 119, Sep. 2020.
- [13] N. Tahir, C.-R. Jung, S.-D. Lee, N. Azizah, W.-C. Ho, and T.-C. Li, "Federated Learning-Based Model for Predicting Mortality: Systematic Review and Meta-Analysis," *J. Med. Internet Res.*, vol. 27, e65708, Jul. 2025. doi: 10.2196/65708
- [14] Xu J, Glicksberg BS, Su C, Walker P, Bian J, Wang F. "Federated learning for healthcare informatics." *Journal of Healthcare Informatics Research*. 2021;5(1):1-19. <https://doi.org/10.1007/s41666-020-00082-4>
- [15] D. C. Nguyen, Q.-V. Pham, P. N. Pathirana, M. Ding, A. Seneviratne, Z. Lin, O. Dobre, and W. Hwang, 'Federated learning for smart healthcare: A survey,' *ACM Computing Surveys*, vol. 55, no. 3, pp. 1–37, Apr. 2023.
- [16] S. Caton and C. Haas, "Fairness in Machine Learning: A Survey," *ACM Comput. Surv.*, vol. 56, no. 7, Art. no. 166, pp. 1–38, Apr. 2024. doi: 10.1145/3616865
- [17] T. Li, M. Sanjabi, A. Beirami, and V. Smith, 'Fair resource allocation in federated learning,' in *Proc. 8th Int. Conf. Learning Representations (ICLR)*, Addis Ababa, Apr. 2020.

- [18] A. Papadaki, N. Martinez, M. Bertran, G. Sapiro, and M. Rodrigues, 'Minimax demographic group fairness in federated learning,' in Proc. 2022 ACM Conf. Fairness, Accountability, and Transparency (FAccT), Seoul, Jun. 2022, pp. 142–159.
- [19] W. Du, T. Xu, X. Wu, and H. Tong, 'Fairness-aware agnostic federated learning,' in Proc. SIAM Int. Conf. Data Mining (SDM), Minneapolis, Apr. 2021, pp. 181–189.
- [20] Y. H. Ezzeldin, S. Yan, C. He, E. Ferrara, and A. S. Avestimehr, 'FairFed: Enabling group fairness in federated learning,' in Proc. 37th AAAI Conf. Artificial Intelligence, Washington D.C., Feb. 2023, pp. 7494–7502.
- [21] M. Abadi et al., 'Deep learning with differential privacy,' in Proc. 2016 ACM SIGSAC Conf. Computer and Communications Security (CCS), Vienna, Oct. 2016, pp. 308–318.
- [22] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang, 'Learning differentially private recurrent language models,' in Proc. 6th Int. Conf. Learning Representations (ICLR), Vancouver, May 2018.
- [23] N. Agarwal, P. Kairouz, and Z. Liu, 'The Skellam mechanism for differentially private federated learning,' in Proc. 35th Conf. Neural Information Processing Systems (NeurIPS), Online, Dec. 2021.
- [24] A. Girgis, D. Data, S. Diggavi, P. Kairouz, and A. T. Suresh, 'Shuffled model of differential privacy in federated learning,' in Proc. 24th Int. Conf. Artificial Intelligence and Statistics (AISTATS), Online, Apr. 2021, pp. 2521–2529.
- [25] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Smola, and V. Smith, 'Federated optimization in heterogeneous networks,' in Proc. Machine Learning and Systems (MLSys), Austin, TX, Mar. 2020, pp. 429–450.
- [26] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, 'SCAFFOLD: Stochastic controlled averaging for federated learning,' in Proc. 37th Int. Conf. Machine Learning (ICML), Online, Jul. 2020, pp. 5132–5143.



AL- Rafidain University

PISSN: (1681-6870); EISSN: (2790-2293)

**Journal of AL-Rafidain**  
**University College for Sciences**

Available online at: <https://www.jruc.s.iq>**JRUCS**Journal of AL-Rafidain  
University College  
for Sciences

## تنقيب البيانات الاتحادي العادل والحافظ للخصوصية للتنبؤ المبكر بالأزمات الصحية: دراسة تجريبية على مستشفيات بغداد

ميادة جبار كيلان

[mayadajabbar@uomustansiriyah.edu.iq](mailto:mayadajabbar@uomustansiriyah.edu.iq)

قسم الدراسات الفندقية، كلية العلوم السياحية، الجامعة المستنصرية، بغداد، العراق.

**المستخلص**

تعد الأزمات الصحية المبكرة، مثل بدء الإنفلونزا، وتدهور الحالة القلبية الوعائية، والفشل التنفسي، مشكلة كبيرة ولم تُحل إلى حد بعيد في الأنظمة الصحية ذات الموارد المحدودة، ولا سيما في البلدان النامية مثل العراق. وتتطلب إجراءات تنقيب البيانات المركزية التقليدية جمع البيانات من جميع المستخدمين وتجميعها في مكان واحد، مما قد يؤدي إلى نتائج متحيزة في تنقيب البيانات تُلحق الضرر بالفئات السكانية قليلة التمثيل. في هذه الورقة، نقدم FedFairMine، وهو نموذج اتحادي لتنقيب البيانات يحافظ على الخصوصية ويعزز العدالة، ويهدف إلى تحسين التنبؤ المبكر بالأزمات الصحية في المستشفيات. ويتضمن FedFairMine ثلاثة ابتكارات رئيسية لتحسين كفاءته (i) محركًا محليًا لاستخراج السمات يولد تمثيلات قائمة على التدرج وذات صلة سريرية مباشرة على الجهاز، بحيث لا تُرسل بيانات المرضى الخام إلى الخادم؛ (ii) آلية تجميع اتحاديّة موجهة نحو العدالة، تهدف إلى تقليل خسارة التنبؤ عبر بيانات البيانات المنفصلة المختلفة، وتحقيق تكافؤ عادل في النتائج بين الفئات الديموغرافية المختلفة، مثل العمر والجنس والوضع الاجتماعي الاقتصادي؛ (iii) نظامًا أنيقًا للخصوصية التفاضلية صُمم خصيصًا للبيانات الحساسة الخاصة بالسجلات الصحية الإلكترونية في المستشفيات العراقية. وقد جرى تقييم FedFairMine على شبكة اتحاديّة تضم ستة من كبرى المستشفيات التعليمية في بغداد: الكندي، وغازي الحريري، وابن سينا، ومدينة الطب، واليرموك، والكاظمية. وفي هذه الدراسة، حللنا 47,312 حالة مرضية منزوعة الهوية بين عامي 2021 و2024. وحقق FedFairMine درجة F1 ممتازة بلغت 0.912، متفوقًا على FedAvg و FedProx و Scaffold بنسبة 8.3% و 5.1% و 3.7% على التوالي. كما خفّض التفاوت الديموغرافي (Equalized Odds difference) إلى 0.021، وهو أفضل بنسبة 76% من معايير تنقيب البيانات المركزية المرجعية. وتؤكد هذه النتائج الإمكانيات الكبيرة لتنقيب البيانات الاتحادي في تحقيق الدقة التنبؤية والعدالة مع مراعاة خصوصية المرضى، وهي خطوة مهمة نحو تحقيق هدفنا الأمم المتحدة للتنمية المستدامة 3 و 10.

**معلومات البحث****تواريخ البحث:**

تاريخ تقديم البحث: 2026/4/20

تاريخ قبول البحث: 2026/5/17

تاريخ رفع البحث على الموقع: 2026/6/30

**الكلمات المفتاحية:**

التعلم الاتحادي، تنقيب البيانات، التنبؤ بالأزمات الصحية، العدالة الخوارزمية، الخصوصية التفاضلية، السجلات الصحية الإلكترونية، مستشفيات بغداد، الحوسبة الطرفية، تكافؤ الفرص، FedFairMine.

**للمراسلة:**

ميادة جبار كيلان

[mayadajabbar@uomustansiriyah.edu.iq](mailto:mayadajabbar@uomustansiriyah.edu.iq)DOI: <https://doi.org/10.55562/jruc.s.v59i1.19>