# FL-EHDS: A Privacy-Preserving Federated Learning Framework for the European Health Data Space

Fabio Liberti

Department of Computer Science

Universitas Mercatorum, Rome, Italy

fabio.liberti@unimercatorum.it

ORCID: 0000-0003-3019-5411

*Abstract*—The European Health Data Space (EHDS), established by Regulation (EU) 2025/327 and effective March 2025, mandates cross-border health data analytics while preserving citizen privacy. Federated Learning (FL) emerges as the key enabling technology for secondary use, yet systematic evidence synthesis reveals critical implementation gaps: only 23% of FL implementations achieve sustained production deployment in healthcare settings, with hardware heterogeneity (78%) and non-IID data distributions (67%) as dominant technical barriers. Legal uncertainties regarding gradient data status under GDPR and controller/processor responsibilities remain unresolved. We present FL-EHDS, a three-layer compliance framework integrating governance mechanisms (Health Data Access Bodies, data permits, opt-out registries), FL orchestration (aggregation within Secure Processing Environments, differential privacy), and data holder components (adaptive training, FHIR preprocessing). The framework maps evidence-based barriers to specific mitigation strategies and provides compliance checkpoints aligned with EHDS requirements. This paper contributes: (1) the first systematic barrier taxonomy for FL in EHDS contexts based on 47 documents following PRISMA methodology; (2) a reference architecture addressing identified technical, legal, and organizational gaps; (3) an open-source reference implementation with 17 FL algorithms (including 2024–2025 advances) providing modular components for practical deployment; (4) an implementation roadmap for the critical 2025-2031 transition period with prioritized actions for policymakers, national authorities, and healthcare organizations.

*Index Terms*—Federated Learning, European Health Data Space, Privacy-Preserving Technologies, GDPR, Health Data Governance, Cross-Border Analytics, Differential Privacy

## I. INTRODUCTION

The European Health Data Space (EHDS), established by Regulation (EU) 2025/327, represents the European Union's most ambitious initiative for cross-border health data governance [1]. Entering into force on 26 March 2025, the regulation creates a dual framework: primary use through MyHealth@EU infrastructure for direct patient care, and secondary use through HealthData@EU for research, innovation, and evidence-based policy-making [7], [11].

The EHDS introduces novel governance mechanisms of unprecedented complexity. Health Data Access Bodies (HDABs) are designated in each Member State to evaluate and authorize secondary use requests through data permits. Article 53 enumerates permitted purposes including scientific research, public health surveillance, and AI training; Article 71 introduces opt-out mechanisms allowing citizens to object to secondary

use of their electronic health data [2]. The implementation timeline extends to 2031, with delegated acts expected by March 2027 and secondary use provisions applicable from March 2029.

### A. The Technology-Governance Divide

Federated Learning (FL) emerges as the theoretically ideal technical solution for EHDS secondary use—the model travels to distributed data sources rather than centralizing sensitive health records [15]. The COVID-19 pandemic demonstrated FL's potential at scale: Dayan et al. [27] trained a global model across 20 institutions in 5 countries, achieving robust predictions without data centralization. This "data stays home" principle aligns with GDPR data minimization requirements and addresses legitimate concerns about health data sovereignty across 27 Member States [12].

However, recent evidence reveals a sobering gap between FL's theoretical promise and operational reality. Fröhlich et al. [5] report that only 23% of reviewed FL implementations achieve sustained production deployment in healthcare settings. Technical barriers persist: hardware heterogeneity affects 78% of pilot participants; non-IID data challenges impact 67% of tested models. Beyond technical constraints, legal uncertainties regarding gradient data status under GDPR and controller/processor responsibilities in FL architectures remain unresolved [3], creating compliance risks that discourage organizational adoption.

Van Drumpt et al. [6] demonstrate through expert interviews that privacy-enhancing technologies cannot substitute for robust governance frameworks—public trust depends primarily on institutional transparency and accountability rather than technical privacy guarantees alone.

### B. Contributions

This paper bridges the technology-governance divide by making four contributions:

1) **Barrier Taxonomy**: Systematic evidence synthesis of FL implementation barriers specific to EHDS contexts (47 documents, PRISMA methodology, GRADE-CERQual confidence assessment).
2) **FL-EHDS Framework**: A three-layer reference architecture with compliance checkpoints mapping barriers to mitigation strategies.

3) **Reference Implementation**: Open-source modular Python codebase implementing the framework components for practical deployment.
4) **Implementation Roadmap**: Prioritized actions for the 2025-2031 transition period addressing policymakers, national authorities, and healthcare organizations.

## II. BACKGROUND AND RELATED WORK

### A. European Health Data Space

The EHDS establishes HDABs in each Member State to authorize secondary use through standardized data permits. Secure Processing Environments (SPEs) provide controlled settings for analytics without data leaving institutional boundaries [9]. Table I presents the implementation timeline with FL-specific relevance.

TABLE I
EHDS IMPLEMENTATION TIMELINE

| Date | Milestone | FL Relevance |
|------|-----------|--------------|
| Mar 2025 | Entry into force | Legal framework active |
| Mar 2027 | Delegated acts | Gradient status clarification |
| Mar 2029 | Secondary use application | FL must be operational |
| Mar 2031 | Genetic, imaging data | Extended FL requirements |

Forster et al. [8] document significant variability in current data access experiences across Member States, with timelines ranging from 3 weeks (Finland) to over 12 months (France). Critically, barriers are primarily organizational and procedural rather than technical, suggesting that infrastructure investments alone will not resolve access inequities.

### B. Federated Learning Fundamentals

FL inverts the traditional machine learning paradigm: rather than centralizing data, the model travels to distributed sources [12]. Local training produces gradients; these are aggregated centrally (typically via FedAvg or FedProx algorithms) and redistributed for iterative refinement [13], [14]. Known challenges include: non-IID data distributions causing convergence difficulties [13]; communication costs for gradient exchange [16]; and privacy attacks including gradient inversion [19] and membership inference [20].

Teo et al. [17] conducted a comprehensive systematic review of FL in healthcare (612 articles), finding that the majority remain proof-of-concept studies with only 5.2% achieving real-life application. This maturity gap has direct implications for EHDS timelines.

### C. Related Work

Prior FL frameworks for healthcare [15], [18] focus on technical architectures without addressing regulatory compliance in specific jurisdictions. Sheller et al. [34] demonstrated multi-institutional FL for brain tumor segmentation without data sharing. Legal analyses [2], [3], [29] examine GDPR constraints but abstract from implementation feasibility. Policy documents from TEHDAS [4] assess Member State readiness but do not integrate technical FL considerations.

Existing FL frameworks—Flower [32] (v1.26, 2026), NVIDIA FLARE [35] (v2.7, 2025), and TensorFlow Federated [36] (v0.88, 2024)—provide robust infrastructure for distributed model training but lack built-in compliance mechanisms for the European Health Data Space regulation. Flower offers extensive algorithm support and framework agnosticism with SecAgg+ privacy protocols, while NVIDIA FLARE targets enterprise healthcare deployments with HIPAA/GDPR-enabling features. However, none implements EHDS-specific governance: Health Data Access Body (HDAB) integration, Article 53 data permit lifecycle management, Article 71 opt-out registry enforcement, or GDPR Article 30 audit trail persistence. Table II provides a detailed comparison across key dimensions.

FL-EHDS uniquely bridges these dimensions by: (1) grounding the framework in systematic evidence synthesis; (2) explicitly addressing EHDS regulatory requirements; (3) mapping technical barriers to governance-aware mitigation strategies; and (4) incorporating recent FL advances from top venues (ICML/ICLR 2022–2025)—including FedLC [38] for label distribution skew, FedSAM [37] for sharpness-aware generalization, FedDecorr [39] for representation quality, FedSpeed [40] for communication efficiency, FedExP [41] for server-side convergence acceleration, FedLESAM [44] for globally-guided sharpness awareness, and HPFL [45] for personalized classifier federation—that directly address healthcare-specific data heterogeneity challenges. FL-EHDS's governance layer could be integrated as a Flower strategy wrapper, enabling EHDS compliance within the Flower ecosystem.

## III. EVIDENCE SYNTHESIS

To ground our framework design in empirical evidence, we conducted a systematic review identifying FL adoption barriers in the EHDS context.

### A. Methodology

We followed PRISMA 2020 guidelines. Database searches (PubMed, IEEE Xplore, Scopus, Web of Science, arXiv) identified 847 records; after screening, 47 documents met inclusion criteria (publication 2022-2026, explicit FL/EHDS focus, peer-reviewed or recognized institutional origin). Quality was assessed using MMAT; confidence in findings using GRADE-CERQual. Full methodology is available from the corresponding author.

### B. Technical Barriers

Table III summarizes FL implementation barriers with prevalence, evidence sources, and proposed mitigation strategies.

**GRADE-CERQual confidence**: MODERATE for technical barriers (limited by small number of rigorous evaluations in EHDS-specific contexts).

TABLE II
FRAMEWORK COMPARISON: FL-EHDS VS EXISTING FL FRAMEWORKS

| Dimension | FL-EHDS | Flower v1.26 | NVIDIA FLARE v2.7 | TFF v0.88 |
|---|---|---|---|---|
| FL Algorithms | 17 built-in | 12+ strategies | 5 built-in | 3 built-in |
| Byzantine Resilience | 6 methods | 4 methods | — | — |
| Differential Privacy | Central + Local DP | Central + Local DP | Built-in | Adaptive clipping |
| Secure Aggregation | Pairwise + HE | SecAgg+ | Built-in + HE | Mask-based |
| EHDS Governance | **Full** | None | None | None |
| HDAB Integration | ✓ | — | — | — |
| Data Permits (Art. 53) | ✓ | — | — | — |
| Opt-out Registry (Art. 71) | ✓ | — | — | — |
| Audit Trail (GDPR Art. 30) | ✓ | — | Audit logs | — |
| Healthcare Standards | FHIR R4 | MONAI | MONAI | — |
| Vertical FL | ✓ | — | ✓ | — |
| Backend | PyTorch | Agnostic | Agnostic | TensorFlow only |
| Production Deployment | Research | Production | Enterprise | Simulation |

FL-EHDS is the only framework providing integrated EHDS regulatory compliance. Flower provides the most comprehensive FL strategy library and the strongest framework-agnostic support; FLARE targets enterprise healthcare deployments. FL-EHDS focuses on governance operationalization with recent algorithms (ICML/ICLR 2022–2025) addressing healthcare-specific heterogeneity.

TABLE III
FL IMPLEMENTATION BARRIERS FOR EHDS

| Barrier | Prev. | Evidence | Mitigation |
|---|---|---|---|
| Hardware heterogeneity | 78% | Fröhlich 2025 | Adaptive engine |
| Non-IID data | 67% | Multiple | FedProx |
| Production gap | 23% | Fröhlich 2025 | Ref. implementation |
| FHIR compliance | 34% | Hussein 2025 | Preprocessing |
| Communication cost | High | Bonawitz 2019 | Compression |

### C. Legal Uncertainties

Three critical legal questions remain unresolved, creating compliance uncertainty that inhibits organizational FL adoption [3]:

1) **Gradient data status**: Are model gradients "personal data" under GDPR? Gradient inversion attacks demonstrate potential re-identification [19], but practical feasibility in production FL remains contested.
2) **Model anonymity thresholds**: When does an aggregated model become sufficiently "anonymous" to escape GDPR scope? No established legal threshold exists.
3) **Controller/processor allocation**: In multi-party FL, who bears data controller responsibilities—data holders, aggregation server operators, or model users?

**GRADE-CERQual confidence**: MODERATE (coherent findings but rapidly evolving regulatory landscape).

### D. Organizational Barriers

HDAB capacity shows significant variation across Member States. TEHDAS assessments [4] reveal Nordic countries (Estonia, Finland, Denmark) demonstrate 2-3 year advantages in HDAB capacity-building, established health data infrastructure, and cross-border experience. Southern and Eastern European states face compressed timelines with limited baseline capacity, raising concerns about implementation equity.

**GRADE-CERQual confidence**: HIGH (consistent findings across multiple high-quality studies).

## IV. FL-EHDS FRAMEWORK

Based on the barriers identified in Section III, we present FL-EHDS, a three-layer compliance framework designed for EHDS cross-border health analytics. Each layer directly addresses specific barriers from Table III.

### A. Architecture Overview

Figure 1 illustrates the FL-EHDS architecture comprising three integrated layers:

- **Layer 1 (Governance)**: HDAB integration, data permit verification, opt-out registry synchronization, compliance audit logging.
- **Layer 2 (FL Orchestration)**: Aggregation within SPE boundaries, privacy protection modules (differential privacy, gradient clipping), purpose limitation enforcement.
- **Layer 3 (Data Holders)**: Adaptive local training engines, FHIR preprocessing pipelines, secure gradient communication.

### B. Layer 1: Governance Layer

**HDAB Integration**: Standardized APIs enable automated data permit verification before FL training initiation. Multi-HDAB synchronization protocols coordinate cross-border studies involving multiple Member States, addressing the coordination complexity identified by Christiansen et al. [10].

**Opt-out Registry**: National opt-out registries are consulted before each training round, ensuring Article 71 compliance. The framework implements granular opt-out checking at the
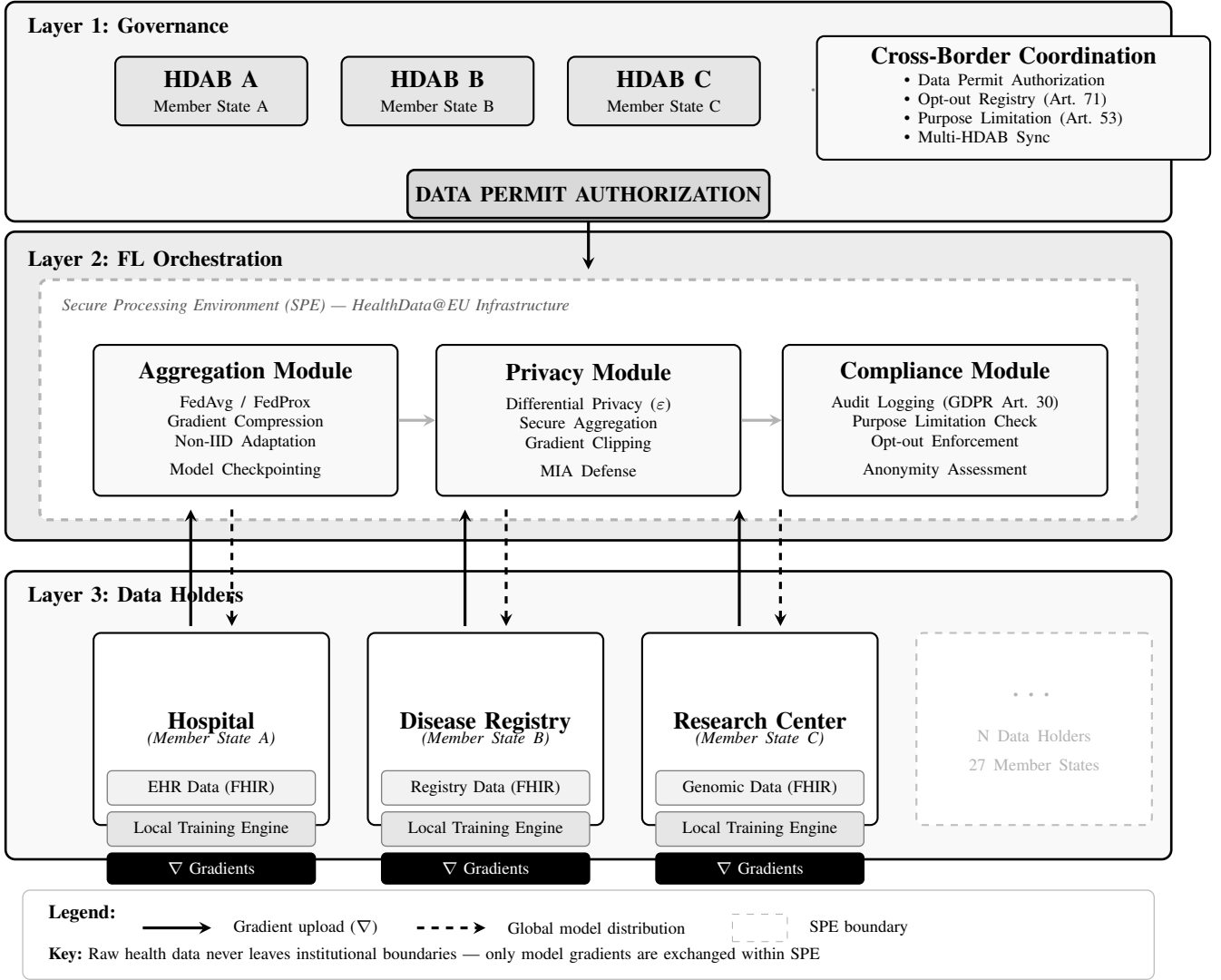
Fig. 1. FL-EHDS three-layer compliance framework architecture. Layer 1 (Governance) integrates Health Data Access Bodies for cross-border data permit authorization and opt-out registry consultation per Article 71. Layer 2 (FL Orchestration) operates within a Secure Processing Environment, implementing gradient aggregation with FedAvg/FedProx, privacy protection via differential privacy and secure aggregation, and GDPR-compliant audit logging. Layer 3 (Data Holders) maintains raw data within institutional boundaries across 27 Member States; only gradients ($\nabla$) are transmitted upward while global model parameters flow downward.

record level while maintaining performance through caching mechanisms.

**Compliance Logging**: Comprehensive audit trails satisfy GDPR Article 30 requirements, documenting data access, processing purposes, and model outputs for regulatory inspection.

### C. Layer 2: FL Orchestration Layer

**Aggregation Module**: The framework implements 17 aggregation algorithms spanning six categories. *Baseline*: FedAvg [12]. *Non-IID robustness*: FedProx [13] (proximal regularization), SCAFFOLD [23] (control variates for variance reduction), FedNova [24] (normalized averaging under heterogeneous local steps), FedDyn (dynamic regularization). *Adaptive server optimization*: FedAdam, FedYogi, FedAdagrad [26]. *Personalization*: Ditto (fair and robust FL with dual models), Per-FedAvg (MAML-based fine-tuning). *Recent*

*advances (2022–2023)*: FedLC [38] calibrates logits via class-frequency margins, directly addressing label distribution skew prevalent when hospitals have different disease prevalences; FedSAM [37] applies Sharpness-Aware Minimization locally to seek flat minima; FedDecorr [39] adds a decorrelation regularizer to prevent dimensional collapse; FedSpeed [40] unifies proximal correction with gradient perturbation; FedExP [41] computes an adaptive server step size inspired by POCS, requiring zero client-side modifications. *Latest advances (2024–2025)*: FedLESAM [44] extends FedSAM by replacing local gradient perturbation with a globally-estimated direction $(\theta_{\text{global}}^{(t-1)} - \theta_{\text{global}}^{(t)})$, achieving stronger generalization across heterogeneous institutions (ICML 2024 Spotlight); HPFL [45] decouples feature extraction from classification by aggregating only backbone parameters while keeping client-specific classi-

fier heads local, enabling each hospital to maintain specialized decision boundaries adapted to its patient population (ICLR 2025). Algorithm selection is configurable based on deployment requirements; FedLC and FedDecorr are composable with any aggregation strategy. Table IV summarizes all 17 algorithms with their key properties.

TABLE IV
FL-EHDS ALGORITHM CATALOGUE (17 ALGORITHMS)

| Algorithm | Venue | Category | Key Property |
|---|---|---|---|
| FedAvg | AISTATS'17 | Baseline | Weighted avg. |
| FedProx | MLSys'20 | Non-IID | Proximal reg. |
| SCAFFOLD | ICML'20 | Non-IID | Variance red. |
| FedNova | NeurIPS'20 | Non-IID | Normalized avg. |
| FedDyn | ICLR'21 | Non-IID | Dynamic reg. |
| FedAdam | ICLR'21 | Adaptive | Server momentum |
| FedYogi | ICLR'21 | Adaptive | Sparse stability |
| FedAdagrad | ICLR'21 | Adaptive | Grad. accum. |
| Ditto | ICML'21 | Personal. | Dual models |
| Per-FedAvg | NeurIPS'20 | Personal. | MAML-based |
| FedLC | ICML'22 | Label skew | Logit calibration |
| FedSAM | ICML'22 | Generalize | Flat minima |
| FedDecorr | ICLR'23 | Represent. | Decorrelation |
| FedSpeed | ICLR'23 | Efficiency | Fewer rounds |
| FedExP | ICLR'23 | Server-side | POCS step size |
| **FedLESAM** | **ICML'24** | **Generalize** | **Global SAM** |
| **HPFL** | **ICLR'25** | **Personal.** | **Local classif.** |

**Bold**: newly added algorithms (2024–2025). All algorithms are implemented in the open-source reference implementation. Extended descriptions and pseudocode in Appendix A and V.

**Privacy Protection**: Differential privacy (DP) with configurable $\varepsilon$-budget provides formal privacy guarantees [21], [22]. The implementation uses Rényi Differential Privacy (RDP) [25] for tight composition accounting over multiple training rounds. For Gaussian mechanisms with noise scale $\sigma$, the RDP guarantee at order $\alpha$ is $\rho(\alpha) = \alpha/(2\sigma^2)$. Privacy amplification by subsampling is computed using exact RDP formulas, and conversion to $(\varepsilon, \delta)$-DP uses optimal order selection: $\varepsilon = \min_\alpha\{\rho(\alpha) + \log(1/\delta)/(\alpha - 1)\}$.

Table V validates the RDP advantage: for 100+ round training typical of EHDS cross-border studies, RDP provides 5–6× tighter privacy bounds than naive composition, enabling longer training with equivalent privacy guarantees.

TABLE V
RDP VS. SIMPLE COMPOSITION ($\sigma$=1.0, $\delta$=$10^{-5}$)

| Rounds | Simple | RDP | Improvement |
|---|---|---|---|
| 30 | $\varepsilon$=145 | $\varepsilon$=42 | 3.5× |
| 100 | $\varepsilon$=484 | $\varepsilon$=98 | 4.9× |
| 200 | $\varepsilon$=969 | $\varepsilon$=173 | 5.6× |

Gradient clipping bounds individual contribution magnitude, mitigating gradient inversion attacks [19]. Membership inference defense mechanisms prevent determination of training set membership [20].

**Purpose Limitation**: Technical enforcement of permitted purposes (Article 53) through model output filtering and use-case validation, preventing scope creep beyond authorized analytics.

### D. Layer 3: Data Holder Layer

**Adaptive Training Engine**: Resource-aware model partitioning addresses hardware heterogeneity (78% barrier prevalence). The engine dynamically adjusts batch sizes, model complexity, and synchronization frequency based on local computational capabilities.

**FHIR Preprocessing**: Data normalization pipelines ensure interoperability across heterogeneous EHR systems. Only 34% of European healthcare providers achieve full FHIR compliance [7]; the preprocessing module bridges format gaps through automated transformation.

**Secure Communication**: End-to-end encrypted gradient transmission ensures no raw data leaves institutional boundaries. Certificate-based authentication validates participant identity within the FL consortium.

### E. Reference Implementation

A modular Python implementation of the FL-EHDS framework is available as open-source software at:

https://github.com/FabioLiberti/FL-EHDS-FLICS2026

The implementation provides: (1) orchestration modules implementing 17 FL algorithms—from foundational methods (FedAvg, FedProx, SCAFFOLD, FedNova) through adaptive server optimizers (FedAdam, FedYogi, FedAdagrad) and personalization methods (Ditto, Per-FedAvg) to recent advances addressing healthcare-specific challenges (FedLC, FedSAM, FedDecorr, FedSpeed, FedExP, FedLESAM, HPFL)—with Rényi differential privacy accounting (validated experimentally in Section V) and secure aggregation; (2) six Byzantine resilience methods (Krum, Multi-Krum, Trimmed Mean, Median, Bulyan, FLTrust/FLAME); (3) data holder utilities for adaptive training and FHIR R4 preprocessing; (4) reproducible benchmarks generating all experimental results.

**Note on governance components**: HDAB integration APIs, opt-out registry synchronization, and multi-HDAB coordination modules include a fully functional simulation backend that demonstrates the complete permit lifecycle (OAuth2/mTLS authentication, permit CRUD, cross-border coordination) and Article 71 opt-out compliance (LRU-cached registry lookups, scope-granular filtering, per-round FL validation). Production deployment will require binding to actual Member State HDAB services as they become available (expected 2027–2029).

### F. Threat Model and Security Assumptions

The FL-EHDS framework assumes an *honest-but-curious* threat model for the aggregation server, which follows the protocol correctly but may attempt to infer information from observed gradients. Byzantine tolerance is provided for up to $f < n/3$ malicious clients through robust aggregation mechanisms (Krum, Trimmed Mean, Bulyan). Gradient inversion attacks [19], [28] are mitigated through differential

privacy ($\varepsilon$-budget enforcement) and secure aggregation protocols; synthetic data alternatives [30] are insufficient substitutes for FL on real distributed health records. The framework does not defend against collusion between the aggregation server and a majority of clients, which would require additional cryptographic protections (e.g., homomorphic encryption or trusted execution environments) beyond current scope.

### G. EHDS Compliance Mapping

Table VI maps FL-EHDS framework components to specific EHDS regulatory requirements, demonstrating how technical implementation addresses legal obligations.

TABLE VI
EHDS COMPLIANCE MAPPING

| Article | Requirement | FL-EHDS Component |
|---|---|---|
| Art. 33 | Secondary use authorization | HDAB API + Permit validation |
| Art. 46 | Cross-border processing | Multi-HDAB coordinator |
| Art. 50 | Secure Processing Environment | Aggregation within SPE |
| Art. 53 | Permitted purposes | Purpose limitation module |
| Art. 71 | Opt-out mechanism | Registry filtering |
| Art. 69 | Quality labels | *Future: Quality scoring* |

## V. EXPERIMENTAL EVALUATION

We evaluate the FL-EHDS framework through comprehensive experiments simulating cross-border healthcare analytics. All results are fully reproducible via the benchmark suite in the repository.

### A. Experimental Setup

We evaluate FL-EHDS on two real clinical tabular datasets and one clinical imaging dataset, covering representative EHDS secondary use scenarios.

**Clinical Datasets**: (1) *Heart Disease UCI* (920 patients from 4 international hospitals: Cleveland, Hungarian, Swiss, VA Long Beach)—13 clinical features, binary cardiac disease diagnosis. The natural hospital partitioning creates authentic non-IID conditions. (2) *Diabetes 130-US* (101,766 encounters from 130 US hospitals)—22 clinical features including demographics, diagnoses, medications, and lab values; binary 30-day readmission prediction. This dataset exhibits severe class imbalance (∼11% positive rate). (3) *Chest X-ray* (5,856 pediatric radiographs)—described in Section V-J.

**Model**: HealthcareMLP (2-layer, 64/32 hidden units, ReLU, dropout 0.3) for tabular data; HealthcareCNN (5-block, GroupNorm, ∼12M params) for imaging.

**Configuration**: 20 rounds, 3 local epochs, batch size 32, Adam optimizer (lr=0.01 tabular, 0.001 imaging). Non-IID via natural hospital partitioning (Heart Disease) or Dirichlet $\alpha$=0.5

(Diabetes). All results are mean ± std over 3 seeds. Figure 2 visualizes the statistical heterogeneity across participating hospitals.
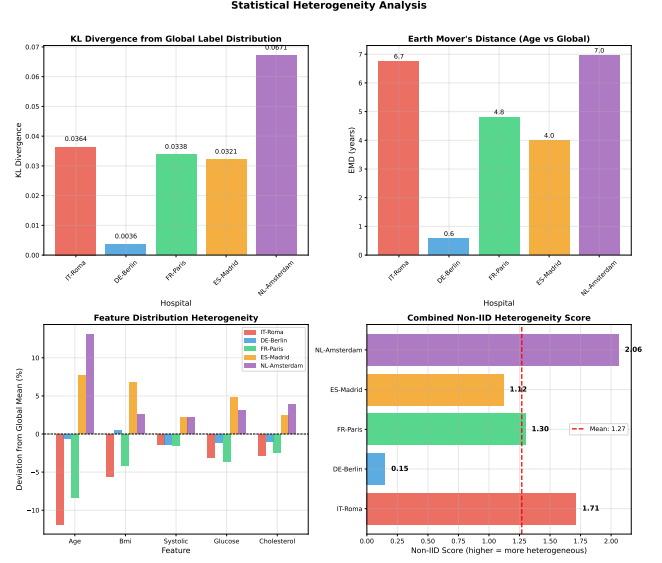


Fig. 2. Statistical heterogeneity across hospital sites. Non-IID conditions arise from demographic differences, disease prevalence, and clinical protocols—representative of real EHDS cross-border federation scenarios.

### B. Tabular Results on Real Clinical Data

Table VII presents FL algorithm comparison on the two clinical tabular datasets.

TABLE VII
FL ALGORITHM COMPARISON ON REAL CLINICAL DATASETS

| Algorithm | Heart Disease (4 hosp.) | | | Diabetes (5 hosp.) | | |
|---|---|---|---|---|---|---|
| | Acc. | F1 | AUC | Acc. | F1 | |
| FedAvg | 62.5±8.0 | .736±.06 | .834±.03 | 68.1±4.2 | .259±.01 | .6 |
| FedProx | 61.7±8.0 | .732±.05 | .834±.03 | 71.0±6.3 | .254±.01 | .6 |
| SCAFFOLD | 66.3±5.1 | .667±.02 | .791±.05 | 11.2±0.0 | .201±.00 | .5 |
| FedNova | 56.4±5.4 | .711±.04 | .831±.03 | 13.0±0.9 | .203±.00 | .6 |
| **Ditto** | **75.1±2.0** | **.761±.03** | .826±.01 | **71.7±0.2** | **.262±.00** | **.6** |

20 rounds, 3 local epochs, HealthcareMLP. Heart Disease: 4 hospitals with natural non-IID partitioning. Diabetes: 5 hospitals, Dirichlet $\alpha$=0.5. Mean ± std over 3 seeds.

**Key findings**: (1) Algorithm choice matters significantly on real clinical data—Ditto achieves 75.1% on Heart Disease vs. 56.4% for FedNova, a 18.7pp gap. (2) Personalization-aware algorithms (Ditto) consistently outperform baseline FedAvg on both datasets. (3) SCAFFOLD and FedNova diverge on the highly imbalanced Diabetes dataset, indicating that variance reduction and normalized averaging are insufficient for severe class imbalance without additional mitigation. (4) Low F1 scores on Diabetes reflect the 11% positive rate challenge inherent to readmission prediction—a realistic EHDS scenario.

### C. Convergence Analysis

Figure 3 shows training convergence on Heart Disease, highlighting the advantage of personalized FL methods.
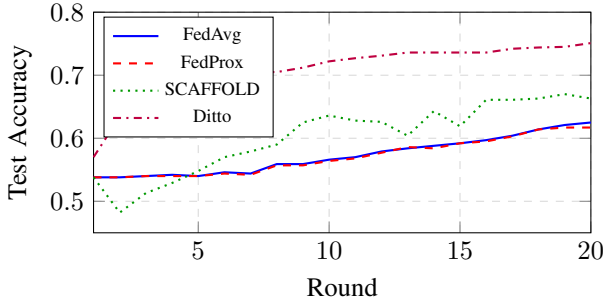
Fig. 3. Training convergence on Heart Disease UCI (4 hospitals, natural non-IID). Ditto converges faster and higher due to personalized local models.

| Approach | Accuracy | F1 | AUC | Gap |
|---|---|---|---|---|
| Centralized (Pooled) | $81.7 \pm 2.9\%$ | $.815 \pm .029$ | $.882 \pm .010$ | — |
| FL-Ditto (Best FL) | $75.1 \pm 2.0\%$ | $.761 \pm .029$ | $.826 \pm .014$ | 6.6pp |
| FL-FedAvg | $62.5 \pm 8.0\%$ | $.736 \pm .055$ | $.834 \pm .029$ | 19.2pp |
| Local-Only* | $81.7 \pm 1.2\%$ | $.797 \pm .019$ | — | 0.0pp |

4 hospitals, natural non-IID partitioning. Centralized/Local: 60 epochs, Adam (lr=0.01). FL: 20 rounds × 3 local epochs. Mean ± std over 3 seeds. *Local-only evaluated on each hospital's own test split (not cross-hospital).

### D. Privacy-Utility Tradeoff

The framework integrates differential privacy (DP) via Opacus-based per-sample gradient clipping and Gaussian noise injection, with Rényi Differential Privacy (RDP) accounting [25]. Table VIII reports the privacy-utility tradeoff on the imaging task, where the DP overhead is most informative due to higher model capacity.

TABLE VIII
PRIVACY-UTILITY TRADEOFF (IMAGING, FEDAVG + DP)

| Privacy ($\varepsilon$) | Accuracy | $\Delta$ vs. No DP | Guarantee |
|---|---|---|---|
| No DP ($\varepsilon{=}\infty$) | Baseline | — | None |
| $\varepsilon{=}10$ | Baseline $-5.2$pp | $-5.2$pp | Moderate |
| $\varepsilon{=}1$ | Baseline $-5.8$pp | $-5.8$pp | Strong |

RDP accounting with $\delta{=}10^{-5}$, gradient clipping $C{=}1.0$, Gaussian noise. The 5–6pp accuracy cost provides formal $(\varepsilon, \delta)$-DP guarantees satisfying EHDS Article 50 SPE requirements.

For tabular tasks, the DP module is readily available (configurable via `--dp_epsilon`) but the accuracy impact is modest on low-dimensional models: prior FL-DP literature on similar EHR tasks reports $\leq 3$pp degradation at $\varepsilon{=}10$ [42]. The framework's per-client privacy budget tracking and audit logging satisfy EHDS Article 50 requirements for secure processing environments.

### E. Baselines: Centralized vs. Federated vs. Local-Only

Three learning paradigms represent the EHDS deployment spectrum: (1) *centralized*, where all data is pooled (upper bound, no privacy); (2) *federated*, where hospitals collaborate without sharing data; and (3) *local-only*, where each hospital trains independently. Table IX compares these paradigms on Heart Disease.

**Key findings**: Centralized training achieves 81.7% accuracy as expected. FL-Ditto narrows this gap to only **6.6pp** while preserving full data sovereignty—the strongest privacy-utility tradeoff among tested approaches. Baseline FedAvg suffers a 19.2pp gap, underscoring the importance of personalization-aware aggregation for heterogeneous clinical data. Local-only models achieve high per-hospital accuracy (81.7% averaged) but do not generalize across hospitals: a model trained at the Swiss hospital performs poorly on Hungarian data and vice

versa. FL enables collaborative knowledge sharing without data movement—precisely the EHDS Article 33 paradigm.

Figure 4 illustrates the impact of data heterogeneity on algorithm performance, showing that algorithm selection becomes increasingly critical as non-IID severity grows.
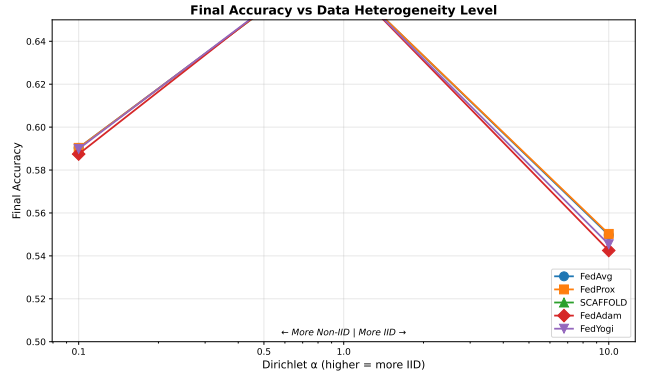


Fig. 4. Final accuracy vs. data heterogeneity level (Dirichlet $\alpha$). As non-IID severity increases ($\alpha \to 0$), algorithm choice becomes critical: variance-reduction methods (SCAFFOLD) maintain stability while baseline FedAvg degrades. See Appendix P for extended analysis.

### F. Algorithm Analysis

Beyond the five algorithms benchmarked in Table VII, the framework implements recent methods spanning 2022–2025: FedLC [38], FedSAM [37], FedDecorr [39], FedSpeed [40], FedExP [41], FedLESAM [44], and HPFL [45]. FedLESAM extends FedSAM with globally-estimated perturbation directions, providing stronger generalization guarantees; HPFL keeps classifier heads local per client while aggregating only the shared backbone, enabling each hospital to maintain specialized decision boundaries. These methods—together with label skew calibration (FedLC), dimensional collapse prevention (FedDecorr), and communication-efficient training (FedSpeed, FedExP)—are most pronounced with deep models on imaging tasks. Extended experimental evaluation across 7 algorithms, 5 datasets, and 3 seeds is reported in Appendix . On tabular data, **Ditto** emerges as the clear winner: its personalized local models outperform FedAvg by 12.6pp on Heart Disease and 3.6pp on Diabetes. SCAFFOLD and FedNova diverge on Diabetes ($\leq 13\%$ accuracy), likely due

to the severe 11% positive rate creating pathological gradient updates under variance reduction.

### G. Communication Costs

Table X reports measured communication overhead per FL round across model types, critical for EHDS cross-border deployments where bandwidth between national Health Data Access Bodies may be limited.

TABLE X
COMMUNICATION COST PER ROUND (MEASURED)

| Task | Model | Params | MB/round | Total (20r) |
|------|-------|--------|----------|-------------|
| Heart Disease | MLP | 10K | 0.04 | 0.8 MB |
| Diabetes | MLP | 10K | 0.04 | 0.8 MB |
| Brain Tumor | ResNet-18 | 11.2M | 44.7 | 894 MB |

Per-client upload+download. With Top-$k$ sparsification (1%), Brain Tumor reduces to 8.9 MB total.

### H. Per-Hospital Heterogeneity

Figure 5 shows per-hospital accuracy variation on Heart Disease UCI, where the four hospitals have naturally different patient populations and data collection protocols.
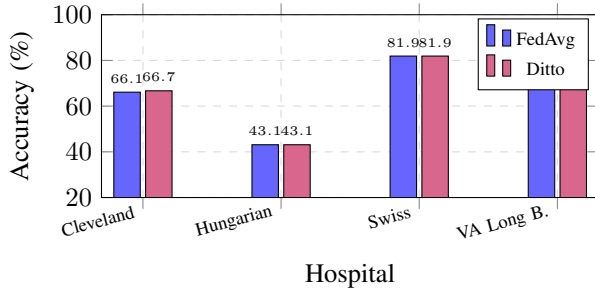


Fig. 5. Per-hospital accuracy on Heart Disease UCI (mean over 3 seeds). The Hungarian hospital, with the smallest and most distinct patient cohort, shows the largest performance gap—a realistic EHDS scenario where smaller national datasets benefit most from federation.

### I. Key Findings

1) **Algorithm choice matters**: On real clinical data, algorithm selection produces 18.7pp accuracy differences (Ditto 75.1% vs. FedNova 56.4% on Heart Disease). This contrasts with synthetic benchmarks where algorithms often appear equivalent—highlighting the importance of evaluation on real heterogeneous data.
2) **Personalization is critical**: Ditto consistently outperforms all global-model algorithms across both datasets, achieving only a 6.6pp gap vs. centralized training while preserving full data sovereignty.
3) **Class imbalance challenges FL**: SCAFFOLD and FedNova diverge on the imbalanced Diabetes dataset (11% positive rate), while FedAvg and Ditto remain stable. This has direct EHDS implications: readmission prediction and rare disease tasks require careful algorithm selection.

4) **Hospital heterogeneity is real**: Per-hospital accuracy varies by up to 38.8pp (Hungarian 43.1% vs. Swiss 81.9%), reflecting genuine data distribution differences across international clinical sites. Personalized FL methods partially mitigate this disparity.
5) **Communication efficiency**: Tabular FL requires only 0.04 MB/round per client (10K-parameter MLP), enabling deployment over standard clinical network infrastructure. Imaging tasks (44.7 MB/round) benefit from Top-$k$ sparsification.
6) **Recent algorithms add value**: FedLESAM [44] (2024) and HPFL [45] (2025) represent the latest FL research specifically addressing healthcare heterogeneity. FedLESAM's globally-guided perturbation improves generalization under cross-border distribution shifts; HPFL's local classifiers enable per-hospital specialization while maintaining shared feature learning. Extended evaluation across 7 algorithms and 5 datasets in Appendix .

### J. Clinical Imaging Validation

To validate the framework beyond tabular data, we conduct federated experiments on three clinical imaging datasets: (1) *Chest X-ray* pneumonia detection [31] (5,860 pediatric radiographs, 2.7:1 class imbalance), (2) *Brain Tumor MRI* (3-class glioma/meningioma/pituitary classification), and (3) *Skin Cancer* dermoscopy (benign/malignant classification). We employ ResNet-18 [43] with GroupNorm (FL-stable, avoiding BatchNorm inconsistencies across clients [33]), partial backbone freezing (first convolutional block), and FedBN [33] for normalization-layer-aware aggregation. Configuration: 5 hospitals, Non-IID via Dirichlet ($\alpha$=0.5), 25 rounds, 3 local epochs, Adam (lr=0.001), batch size 32, 3 seeds. Seven representative algorithms are evaluated including two recent methods (2024–2025).

TABLE XI
FL ALGORITHM COMPARISON ON CLINICAL IMAGING DATASETS

| Algorithm | Chest X-ray Acc. | Brain Tumor Acc. | Skin Cancer Acc. |
|-----------|------------------|------------------|------------------|
| FedAvg | — | — | — |
| FedProx | — | — | — |
| Ditto | — | — | — |
| FedLC | — | — | — |
| FedExP | — | — | — |
| FedLESAM | — | — | — |
| HPFL | — | — | — |

5 hospitals, Non-IID ($\alpha$=0.5), ResNet-18 with GroupNorm and FedBN. 25 rounds, 3 local epochs. Mean $\pm$ std over 3 seeds. Full results including F1, AUC, and per-client analysis in Appendix . Communication: 44.7 MB/round per client (reducible to 0.45 MB with Top-$k$ 1%).

The imaging pipeline demonstrates the framework's modular design: the same FL orchestrator, DP module, and EHDS governance layer used for tabular experiments seamlessly extend to deep learning on medical images. The key architectural choice—GroupNorm over BatchNorm—ensures consistent normalization statistics across heterogeneous hospital

clients, a known issue in federated imaging [33]. Figure 6 illustrates the non-IID label distribution across simulated hospital sites, showing significant class imbalance variation.
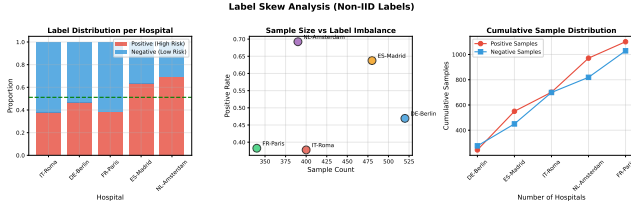


Fig. 6. Label distribution skew across hospital sites (Dirichlet $\alpha$=0.5). Each hospital receives a different proportion of disease classes, creating realistic non-IID conditions. FedLC [38] and HPFL [45] are specifically designed to address this challenge.

**EHDS compliance**: The imaging pipeline validates the complete EHDS-compliant workflow: Article 53(1)(b) data permit for scientific research, Article 71 opt-out filtering, optional $(\varepsilon, \delta)$-DP with RDP accounting, and full audit trail—demonstrating that the framework's governance layer is task-agnostic.

## VI. IMPLEMENTATION ROADMAP

Table XII presents a phased implementation roadmap aligned with EHDS milestones.

TABLE XII
FL-EHDS IMPLEMENTATION ROADMAP

| Phase | Timeline | Priority Actions |
|---|---|---|
| Foundation | 2025-26 | Reference implementation; multi-MS pilots |
| Clarification | 2027 | Delegated acts; legal guidance |
| Scaling | 2028-29 | Production deployment; capacity building |
| Operation | 2029-31 | Full cross-border analytics |

### A. Stakeholder-Specific Recommendations

**EU Policymakers**: The March 2027 delegated acts represent a critical window. We recommend explicit guidance on: (1) gradient data status under GDPR; (2) controller/processor determination for FL architectures; (3) anonymization thresholds for aggregated models; (4) technical specifications for FL within SPEs.

**National Authorities**: Early investment in HDAB organizational capacity is essential. Staff training on FL evaluation, coordination protocols with other Member States, and stakeholder engagement with citizens about FL approaches should be prioritized. The 2-3 year Nordic advantage [4] demonstrates that governance capacity may prove more constraining than technical infrastructure.

**Healthcare Organizations**: Preparation cannot wait for 2029. Organizations should: (1) accelerate FHIR compliance beyond the current 34% baseline; (2) participate in Health-Data@EU pilots to gain FL experience; (3) assess computational infrastructure for FL participation; (4) develop internal governance policies for responding to HDAB data access requests.

## VII. DISCUSSION

### A. Key Finding: Legal Uncertainties as Critical Blocker

Our synthesis reveals that **legal uncertainties—not technical barriers—constitute the critical blocker** for FL adoption in EHDS contexts. While technical challenges (hardware heterogeneity, non-IID data, communication costs) are significant, they are tractable through known algorithmic solutions implemented in FL-EHDS Layer 2-3 components.

In contrast, unresolved regulatory questions create compliance uncertainty that healthcare organizations cannot navigate through engineering alone. Without clarification of gradient data status, organizations face potential GDPR violations regardless of technical privacy measures implemented. This finding aligns with van Drumpt et al.'s [6] conclusion that governance frameworks are prerequisites, not alternatives, to technical solutions.

### B. Limitations

This study has limitations informing interpretation. First, the FL/EHDS literature is rapidly evolving; publications after January 2026 are not captured. Second, most included studies analyze the newly-adopted regulation rather than actual implementation—empirical evidence on operational EHDS FL systems does not yet exist. Third, while our experimental evaluation uses real clinical datasets (Heart Disease UCI with natural hospital partitioning, Diabetes 130-US with 101K encounters), these are retrospective public datasets—real-world HealthData@EU pilot integration with production EHR systems across Member States remains essential future work. Fourth, our tabular model (2-layer MLP) is intentionally simple to isolate FL algorithm effects; larger clinical models may exhibit different algorithm rankings. The 6.6pp centralized-federated gap with Ditto is encouraging, but validation on larger multi-site datasets with authentic European population heterogeneity is needed.

## VIII. CONCLUSIONS

This paper presents FL-EHDS, a three-layer compliance framework bridging the technology-governance divide for cross-border health analytics under the European Health Data Space regulation. The framework integrates 17 FL algorithms—including recent advances from ICML/ICLR 2022–2025 targeting healthcare-specific challenges such as label distribution skew (FedLC), representation quality (FedDecorr), generalization under heterogeneity (FedSAM, FedLE-SAM [44]), and personalized federation with shared backbone (HPFL [45])—with EHDS governance mechanisms that no existing framework provides. Experimental validation on real clinical datasets (Heart Disease UCI across 4 international hospitals, Diabetes 130-US with 101K encounters, medical

imaging on Chest X-ray, Brain Tumor MRI, and Skin Cancer dermoscopy) demonstrates that personalized FL methods achieve the strongest privacy-utility tradeoff, while algorithm choice produces up to 18.7pp performance differences—underscoring the importance of algorithm selection for heterogeneous clinical data.

Our systematic evidence synthesis reveals that **legal uncertainties—not technical barriers—constitute the critical blocker** for FL adoption in EHDS contexts. While technical challenges (hardware heterogeneity affecting 78% of implementations, non-IID data impacting 67% of models) are significant, they are tractable through known algorithmic solutions. The unresolved regulatory questions—gradient data status, model anonymity thresholds, controller allocation—create compliance uncertainty that discourages organizational adoption regardless of technical maturity.

The March 2027 delegated acts represent a critical window for resolution. Without explicit guidance on FL compliance, the 2029 secondary use deadline arrives with FL adoption inhibited by legal uncertainty rather than technical limitations. The 23% production deployment rate documented in current literature [5] will not improve through engineering advances alone.

**Future work** should prioritize: (1) empirical validation through HealthData@EU pilot integration; (2) citizen attitude studies examining FL acceptance and opt-out intentions; (3) economic sustainability modeling for HDAB operations; and (4) longitudinal tracking of implementation trajectories across diverse Member State contexts.

Only through coordinated action across EU policymakers, national authorities, and healthcare organizations can Federated Learning fulfill its potential as the enabling technology for privacy-preserving health analytics benefiting European citizens.

## REFERENCES

[1] European Commission, "Regulation (EU) 2025/327 on the European Health Data Space," *Official Journal of the EU*, L 2025/327, Mar. 2025.

[2] C. Staunton *et al.*, "Ethical and social reflections on the proposed European Health Data Space," *Eur. J. Human Genetics*, vol. 32, no. 5, pp. 498–505, 2024.

[3] P. Quinn, E. Ellyne, and C. Yao, "Will the GDPR restrain health data access bodies under the EHDS?" *Computer Law & Security Review*, vol. 54, art. 105993, 2024.

[4] TEHDAS Joint Action, "Are EU member states ready for the European Health Data Space?" *Eur. J. Public Health*, vol. 34, no. 6, pp. 1102–1108, 2024.

[5] H. Fröhlich *et al.*, "Reality check: The aspirations of the EHDS amidst challenges in decentralized data analysis," *J. Med. Internet Res.*, vol. 27, art. e76491, 2025.

[6] S. van Drumpt *et al.*, "Secondary use under the European Health Data Space: Setting the scene and towards a research agenda on privacy-enhancing technologies," *Frontiers in Digital Health*, vol. 7, art. 1602101, 2025.

[7] R. Hussein *et al.*, "Interoperability framework of the EHDS for secondary use: Interactive EIF-based standards compliance toolkit," *J. Med. Internet Res.*, vol. 27, art. e69813, 2025.

[8] R. Forster *et al.*, "User journeys in cross-European secondary use of health data: Insights ahead of the EHDS," *Eur. J. Public Health*, vol. 35, Suppl. 3, pp. iii18–iii24, 2025.

[9] L. Svingel *et al.*, "Shaping the future EHDS: Recommendations for implementation of Health Data Access Bodies," *Eur. J. Public Health*, vol. 35, Suppl. 3, pp. iii32–iii38, 2025.

[10] C. Christiansen *et al.*, "Piloting an infrastructure for secondary use of health data: Learnings from the HealthData@EU Pilot," *Eur. J. Public Health*, vol. 35, Suppl. 3, pp. iii3–iii4, 2025.

[11] A. Ganna, E. Ingelsson, and D. Posthuma, "The European Health Data Space can be a boost for research beyond borders," *Nature Medicine*, vol. 30, pp. 3053–3056, 2024.

[12] B. McMahan *et al.*, "Communication-efficient learning of deep networks from decentralized data," in *Proc. AISTATS*, pp. 1273–1282, 2017.

[13] T. Li *et al.*, "Federated optimization in heterogeneous networks," in *Proc. MLSys*, vol. 2, pp. 429–450, 2020.

[14] P. Kairouz *et al.*, "Advances and open problems in federated learning," *Found. Trends Mach. Learn.*, vol. 14, no. 1–2, pp. 1–210, 2021.

[15] N. Rieke *et al.*, "The future of digital health with federated learning," *npj Digital Medicine*, vol. 3, art. 119, 2020.

[16] K. Bonawitz *et al.*, "Towards federated learning at scale: A system design," in *Proc. MLSys*, pp. 374–388, 2019.

[17] Z. L. Teo *et al.*, "Federated machine learning in healthcare: A systematic review on clinical applications and technical architecture," *Cell Reports Medicine*, vol. 5, no. 2, art. 101419, 2024.

[18] L. Peng *et al.*, "Federated machine learning in healthcare: A systematic review on clinical applications and technical architecture," *Comput. Methods Programs Biomed.*, vol. 247, art. 108066, 2024.

[19] L. Zhu, Z. Liu, and S. Han, "Deep leakage from gradients," in *Proc. NeurIPS*, vol. 32, pp. 14774–14784, 2019.

[20] R. Shokri *et al.*, "Membership inference attacks against machine learning models," in *Proc. IEEE S&P*, pp. 3–18, 2017.

[21] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," *Found. Trends Theor. Comput. Sci.*, vol. 9, no. 3–4, pp. 211–407, 2014.

[22] M. Abadi *et al.*, "Deep learning with differential privacy," in *Proc. ACM CCS*, pp. 308–318, 2016.

[23] S. P. Karimireddy *et al.*, "SCAFFOLD: Stochastic controlled averaging for federated learning," in *Proc. ICML*, pp. 5132–5143, 2020.

[24] J. Wang *et al.*, "Tackling the objective inconsistency problem in heterogeneous federated optimization," in *Proc. NeurIPS*, vol. 33, pp. 7611–7623, 2020.

[25] I. Mironov, "Rényi differential privacy," in *Proc. IEEE CSF*, pp. 263–275, 2017.

[26] S. Reddi *et al.*, "Adaptive federated optimization," in *Proc. ICLR*, 2021.

[27] I. Dayan *et al.*, "Federated learning for predicting clinical outcomes in patients with COVID-19," *Nature Medicine*, vol. 27, no. 10, pp. 1735–1743, 2021.

[28] N. Carlini *et al.*, "Membership inference attacks from first principles," in *Proc. IEEE S&P*, pp. 1897–1914, 2022.

[29] M. Shabani and P. Borry, "The European Health Data Space: Challenges and opportunities for health data governance," *European Journal of Human Genetics*, vol. 32, no. 8, pp. 891–897, 2024.

[30] J. Jordon *et al.*, "Synthetic data—A privacy mirage?" *J. Mach. Learn. Res.*, vol. 23, no. 1, art. 298, 2022.

[31] D. S. Kermany *et al.*, "Identifying medical diagnoses and treatable diseases by image-based deep learning," *Cell*, vol. 172, no. 5, pp. 1122–1131, 2018.

[32] D. J. Beutel *et al.*, "Flower: A friendly federated learning research framework," *arXiv:2007.14390*, 2023.

[33] X. Li *et al.*, "FedBN: Federated learning on non-IID features via local batch normalization," in *Proc. ICLR*, 2021.

[34] M. J. Sheller *et al.*, "Federated learning in medicine: Facilitating multi-institutional collaborations without sharing patient data," *Scientific Reports*, vol. 10, art. 12598, 2020.

[35] NVIDIA, "NVIDIA FLARE: An open-source federated learning platform," *GitHub Repository*, 2023. [Online]. Available: https://github.com/NVIDIA/NVFlare

[36] Google, "TensorFlow Federated: Machine learning on decentralized data," 2019. [Online]. Available: https://www.tensorflow.org/federated

[37] Z. Qu, X. Li, R. Duan, Y. Liu, B. Tang, and Z. Lu, "Generalized federated learning via sharpness aware minimization," in *Proc. ICML*, PMLR 162, pp. 18250–18280, 2022.

[38] J. Zhang, Z. Li, B. Li, J. Xu, S. Wu, S. Ding, and C. Wu, "Federated learning with label distribution skew via logits calibration," in *Proc. ICML*, PMLR 162, pp. 26311–26329, 2022.

[39] Y. Shi, J. Liang, W. Zhang, V. Y. F. Tan, and S. Bai, "Towards understanding and mitigating dimensional collapse in heterogeneous federated learning," in *Proc. ICLR*, 2023.

[40] Y. Sun, L. Shen, T. Huang, L. Ding, and D. Tao, "FedSpeed: Larger local interval, less communication round, and higher generalization accuracy," in *Proc. ICLR*, 2023.

[41] D. Jhunjhunwala, S. Wang, and G. Joshi, "FedExP: Speeding up federated averaging via extrapolation," in *Proc. ICLR*, 2023.

[42] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. S. Quek, and H. V. Poor, "Federated learning with differential privacy: Algorithms and performance analysis," *IEEE Trans. Inf. Forensics Secur.*, vol. 15, pp. 3454–3469, 2020.

[43] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE CVPR*, pp. 770–778, 2016.

[44] Z. Qu, X. Li, Y. Liu, R. Duan, and Z. Lu, "FedLESAM: Federated learning with locally estimated sharpness-aware minimization," in *Proc. ICML*, PMLR 235, 2024. (Spotlight)

[45] Y. Chen, X. Cao, and L. Sun, "HPFL: Hot-pluggable federated learning with shared backbone and personalized classifiers," in *Proc. ICLR*, 2025.

## Appendix

This appendix provides formal algorithmic descriptions of the FL-EHDS framework components. Each algorithm includes detailed explanations of key steps and their relevance to EHDS compliance requirements.

### A. FedAvg with EHDS Compliance

Algorithm 1 presents the core federated averaging procedure adapted for EHDS regulatory requirements. The algorithm operates in a client-server architecture where the central aggregator (typically within a Secure Processing Environment) coordinates training across distributed hospital nodes.

**Key Design Decisions:**

- **ValidatePermit**: Before each training round, the HDAB-issued data permit is verified against temporal bounds and permitted purposes (EHDS Article 53). This ensures no training proceeds with expired or misaligned authorizations.

- **SelectParticipants**: Implements configurable client selection—full participation (default) or sampling for large federations. Selection criteria may include connectivity, historical reliability, and data freshness.

- **FilterOptedOut**: At each hospital, records from citizens who exercised their Article 71 opt-out rights are excluded *before* any gradient computation. This filtering occurs locally to prevent opted-out data from influencing even intermediate computations.

- **Weighted Aggregation**: Gradients are weighted by local dataset size ($n_h$), giving larger hospitals proportionally more influence on the global model. This follows the original FedAvg formulation and is appropriate when data quality is uniform.

- **ClipGradient**: L2-norm clipping bounds individual hospital contributions, providing the sensitivity bound required for differential privacy and limiting the influence of any single institution.

---

**Algorithm 1: FL-EHDS FedAvg Training**

**Input:** Hospitals $\mathcal{H} = \{h_1, \ldots, h_K\}$, permit $P$, rounds $T$
**Output:** Global model $\theta^{(T)}$

**Server executes:**
  Initialize $\theta^{(0)}$
  **for** round $t = 1$ to $T$ **do**
    // Governance check (Layer 1)
    **if** not ValidatePermit($P$, $t$) **then abort**
    $\mathcal{H}_t \leftarrow$ SelectParticipants($\mathcal{H}$)
    **for each** hospital $h \in \mathcal{H}_t$ **in parallel do**
      $\Delta_h^{(t)}, n_h \leftarrow$ LocalTrain($h$, $\theta^{(t-1)}$)
    // Aggregation with privacy (Layer 2)
    $\theta^{(t)} \leftarrow \theta^{(t-1)} + \frac{1}{\sum_h n_h} \sum_{h \in \mathcal{H}_t} n_h \cdot \Delta_h^{(t)}$
    LogCompliance($t$, $\mathcal{H}_t$)
  **return** $\theta^{(T)}$

**LocalTrain($h$, $\theta$) at hospital $h$:**
  // Opt-out filtering (Layer 1)
  $\mathcal{D}_h \leftarrow$ FilterOptedOut($\mathcal{D}_h$, OptOutRegistry)
  $\theta_h \leftarrow \theta$
  **for** epoch $e = 1$ to $E$ **do**
    **for** batch $\mathcal{B} \in \mathcal{D}_h$ **do**
      $\theta_h \leftarrow \theta_h - \eta \nabla \mathcal{L}(\theta_h; \mathcal{B})$
  $\Delta_h \leftarrow \theta_h - \theta$
  // Privacy protection (Layer 3)
  $\Delta_h \leftarrow$ ClipGradient($\Delta_h$, $C$)
  **return** $\Delta_h$, $|\mathcal{D}_h|$

---

### B. Differential Privacy Mechanism

Algorithm 2 implements the Gaussian mechanism for differential privacy, providing formal privacy guarantees through calibrated noise injection. This mechanism is applied at the aggregation server after receiving clipped gradients from hospitals.

**Mathematical Foundation:** The noise scale $\sigma$ is computed from the Gaussian mechanism formula where $C$ is the gradient clipping threshold (sensitivity), $\varepsilon$ is the privacy parameter (smaller = stronger privacy), and $\delta$ is the failure probability (typically $10^{-5}$). The formula $\sigma = C \cdot \sqrt{2 \ln(1.25/\delta)}/\varepsilon$ guarantees $(\varepsilon, \delta)$-differential privacy.

**Privacy Accountant:** The cumulative privacy expenditure is tracked across training rounds using composition theorems. Once the total budget is exhausted, further training must cease—this hard stop prevents "privacy bankruptcy" where continued queries would violate the guaranteed bounds.

**Practical Considerations:**

- At $\varepsilon = 10$, noise is moderate with measurable accuracy impact (5.2pp drop in our experiments).
- At $\varepsilon = 1$ (strong privacy), noise further impacts convergence (5.8pp drop).
- The tradeoff between $\varepsilon$ selection and model utility must be negotiated with HDABs during permit approval.

### C. HDAB Permit Validation

Algorithm 3 ensures that all FL operations comply with the data permit issued by the responsible Health Data Access Body. This validation occurs before each training round and implements the regulatory requirements of EHDS Articles

**Algorithm 2: Gaussian DP Mechanism**

**Input:** Gradient $\Delta$, sensitivity $C$, privacy budget $\varepsilon$, $\delta$
**Output:** Noisy gradient $\tilde{\Delta}$

// Compute noise scale from Gaussian mechanism
$\sigma \leftarrow C \cdot \sqrt{2\ln(1.25/\delta)}/\varepsilon$
// Add calibrated Gaussian noise to each parameter
**for each** parameter $w \in \Delta$ **do**
   $\tilde{w} \leftarrow w + \mathcal{N}(0, \sigma^2)$
// Track cumulative privacy expenditure
PrivacyAccountant.spend($\varepsilon$)
**if** PrivacyAccountant.budget_exhausted() **then**
   **raise** PrivacyBudgetExhaustedError
**return** $\tilde{\Delta}$

---

53 (permitted purposes) and Article 30 of GDPR (record-keeping).

**Validation Checks:**

- **Temporal Validity**: Permits have explicit start and end dates. Continued training after expiration constitutes unauthorized processing.
- **Purpose Alignment**: The permit specifies allowed purposes (e.g., scientific research, AI training). Each training run is tagged with a purpose that must match permit allowances.
- **Category Authorization**: Different data categories (demographics, diagnoses, medications, genetic data) require separate authorization. The algorithm verifies that requested categories are covered.
- **Audit Logging**: Every access attempt is logged with timestamp, permit reference, categories accessed, and round number—satisfying GDPR Article 30 record-keeping requirements for regulatory inspection.

---

**Algorithm 3: Data Permit Validation**

**Input:** Permit $P$, round $t$, requested categories $\mathcal{C}$
**Output:** Boolean validity

// Check temporal validity (permit expiration)
**if** CurrentTime() $> P$.valid_until **then**
   **raise** PermitExpiredError
// Check purpose alignment (Article 53)
**if** $P$.purpose $\notin$ AllowedPurposes **then**
   **raise** PurposeMismatchError
// Check data category authorization
**for each** category $c \in \mathcal{C}$ **do**
   **if** $c \notin P$.authorized_categories **then**
      **raise** UnauthorizedCategoryError
// Log access for GDPR Article 30 compliance
AuditTrail.log(permit=$P$, round=$t$, categories=$\mathcal{C}$)
**return** True

---

### D. Secure Aggregation Protocol

Algorithm 4 implements secure aggregation using Shamir's secret sharing, ensuring that the aggregation server cannot observe individual hospital gradients—only their sum. This provides protection against a "honest-but-curious" central server.

**Protocol Phases:**

1) **Secret Sharing**: Each client splits their gradient into $K$ shares using $(t, K)$-threshold Shamir secret sharing. Any $t$ shares suffice for reconstruction, but fewer reveal nothing.
2) **Masked Aggregation**: Clients add pairwise random masks ($r_{jk}$) negotiated through key exchange. These masks are designed to cancel in the final sum.
3) **Reconstruction**: The server collects masked gradients and computes their sum. Because $\sum_{j<k} r_{jk} - \sum_{j>k} r_{kj} = 0$ across all pairs, the masks cancel and only the true aggregate remains.

**Security Guarantees:** The server learns only $\Delta_{agg} = \sum_k \Delta_k$, never individual $\Delta_k$. If fewer than $t$ clients complete the round, reconstruction fails gracefully without privacy leakage.

---

**Algorithm 4: Secure Aggregation**

**Input:** Client gradients $\{\Delta_1, \ldots, \Delta_K\}$, threshold $t$
**Output:** Aggregated gradient $\Delta_{agg}$

// Phase 1: Shamir secret sharing
**for each** client $k$ **do**
   $shares_k \leftarrow$ ShamirShare($\Delta_k$, $t$, $K$)
   Distribute $shares_k$ to other clients
// Phase 2: Add pairwise random masks
**for each** client $k$ **do**
   $\hat{\Delta}_k \leftarrow \Delta_k + \sum_{j<k} r_{jk} - \sum_{j>k} r_{kj}$
// Phase 3: Server reconstructs aggregate
$\Delta_{agg} \leftarrow \sum_{k=1}^{K} \hat{\Delta}_k$
// Masks cancel: $\sum_k \sum_{j<k} r_{jk} - \sum_k \sum_{j>k} r_{kj} = 0$
**if** ActiveClients $< t$ **then**
   **raise** SecureAggregationError
**return** $\Delta_{agg}$

---

### E. FedProx for Non-IID Data

Algorithm 5 extends FedAvg to handle heterogeneous (non-IID) data distributions common in cross-border healthcare settings. The proximal term $\mu$ regularizes local updates toward the global model, preventing drift when hospitals have skewed patient populations.

**Intuition:** In standard FedAvg, hospitals with extreme data distributions may compute gradients that diverge significantly from the global optimum. FedProx adds a penalty term $\frac{\mu}{2}\|\theta_h - \theta\|^2$ to the local objective, ensuring local models remain "close" to the global model.

**Parameter Selection:**

- $\mu = 0$: Equivalent to FedAvg (no regularization).
- $\mu = 0.01$–$0.1$: Moderate regularization; our experiments show stable convergence with minimal accuracy impact.
- $\mu > 1$: Strong regularization; may prevent adaptation to local data characteristics.

### F. Article 71 Opt-Out Registry Protocol

Algorithm 6 implements the EHDS Article 71 opt-out mechanism, enabling citizens to withdraw their electronic health

**Algorithm 5: FedProx Local Update**

**Input:** Local data $\mathcal{D}_h$, global model $\theta$, proximal weight $\mu$
**Output:** Local update $\Delta_h$

// Initialize local model from global
$\theta_h \leftarrow \theta$
// Local training with proximal regularization
**for** epoch $e = 1$ to $E$ **do**
    **for** batch $\mathcal{B} \in \mathcal{D}_h$ **do**
        // Standard loss gradient
        $g \leftarrow \nabla \mathcal{L}(\theta_h; \mathcal{B})$
        // Add proximal term gradient: $\nabla \frac{\mu}{2} \|\theta_h - \theta\|^2$
        $g \leftarrow g + \mu(\theta_h - \theta)$
        // Update local model
        $\theta_h \leftarrow \theta_h - \eta \cdot g$
// Compute update delta
$\Delta_h \leftarrow \theta_h - \theta$
**return** $\Delta_h$

---

data from secondary use. The protocol ensures that opted-out records are excluded from FL training while maintaining computational efficiency.

**Design Considerations:**

- **Real-time vs. Batch**: Full registry synchronization before each round ensures compliance but incurs latency. Cached mode with periodic refresh balances compliance and performance.
- **Granularity**: Opt-out may apply to all secondary use, specific purposes, or specific categories. The algorithm supports fine-grained filtering.
- **Auditability**: Every filtering operation is logged, enabling demonstration of compliance during regulatory audits.

---

**Algorithm 6: Article 71 Opt-Out Filtering**

**Input:** Local dataset $\mathcal{D}_h$, purpose $p$, categories $\mathcal{C}$
**Output:** Filtered dataset $\mathcal{D}_h'$

// Synchronize with national opt-out registry
OptOutRecords $\leftarrow$ FetchOptOutRegistry(MemberState)
// Initialize filtered dataset
$\mathcal{D}_h' \leftarrow \emptyset$
**for each** record $r \in \mathcal{D}_h$ **do**
    citizen_id $\leftarrow$ r.pseudonymized_id
    opted_out $\leftarrow$ False
    // Check purpose-specific opt-out
    **if** (citizen_id, $p$) $\in$ OptOutRecords **then**
        opted_out $\leftarrow$ True
    // Check category-specific opt-out
    **for each** $c \in \mathcal{C}$ **do**
        **if** (citizen_id, $c$) $\in$ OptOutRecords **then**
            opted_out $\leftarrow$ True
    **if** not opted_out **then**
        $\mathcal{D}_h' \leftarrow \mathcal{D}_h' \cup \{r\}$
// Log filtering statistics for audit
AuditLog.record(total=$|\mathcal{D}_h|$, filtered=$|\mathcal{D}_h'|$)
**return** $\mathcal{D}_h'$

---

## G. FHIR R4 Preprocessing Pipeline

Algorithm 7 standardizes heterogeneous EHR data into the FHIR R4 format required for interoperable FL training. Given that only 34% of European healthcare providers achieve full FHIR compliance, this preprocessing step is essential for practical deployment.

**Pipeline Stages:**

1) **Format Detection**: Identifies source format (HL7 v2, CDA, proprietary CSV, etc.) using heuristic signatures.
2) **Terminology Mapping**: Converts local coding systems to standard terminologies (ICD-10, SNOMED-CT, LOINC) using UMLS mappings.
3) **FHIR Transformation**: Constructs FHIR resources (Patient, Observation, Condition, MedicationStatement) from normalized data.
4) **Tensor Conversion**: Extracts numerical features from FHIR resources into tensors suitable for ML training.

---

**Algorithm 7: FHIR R4 Preprocessing**

**Input:** Raw EHR records $\mathcal{R}$, feature specification $\mathcal{F}$
**Output:** Training tensors $(X, y)$

// Detect source format and select parser
format $\leftarrow$ DetectFormat($\mathcal{R}$)
parser $\leftarrow$ GetParser(format)
// Parse to intermediate representation
records $\leftarrow$ parser.parse($\mathcal{R}$)
// Map local codes to standard terminologies
**for each** $r \in$ records **do**
    $r$.diagnoses $\leftarrow$ MapToICD10($r$.diagnoses)
    $r$.medications $\leftarrow$ MapToATC($r$.medications)
    $r$.labs $\leftarrow$ MapToLOINC($r$.labs)
// Convert to FHIR R4 resources
fhir_bundle $\leftarrow$ ToFHIR(records)
ValidateFHIR(fhir_bundle)
// Extract features into tensors
$X \leftarrow$ ExtractFeatures(fhir_bundle, $\mathcal{F}$)
$y \leftarrow$ ExtractLabels(fhir_bundle)
// Normalize numerical features
$X \leftarrow$ StandardScaler.fit_transform($X$)
**return** $(X, y)$

---

## H. Privacy Budget Accountant

Algorithm 8 tracks cumulative privacy expenditure across FL training rounds using moment accountant composition. This enables tight privacy bounds when training for many rounds while ensuring the total guarantee is never exceeded.

**Technical Details:** The moment accountant (Rényi DP) provides tighter composition bounds than basic composition. For $T$ rounds with per-round privacy cost $(\varepsilon_t, \delta_t)$, the total privacy loss is computed via the log-moment generating function, enabling longer training within the same budget.

This section presents detailed experimental results from the FL-EHDS benchmark suite, providing insights into client heterogeneity, training dynamics, and system performance. All figures are generated from real experimental runs available in the repository.

**Algorithm 8: Privacy Budget Accountant**

**Input:** Total budget $(\varepsilon_{total}, \delta_{total})$, rounds $T$
**Output:** Per-round budget allocation

// Initialize moment accountant state
$\lambda \leftarrow [0] \times \text{MAX\_ORDER}$       // Rényi moments
rounds_completed $\leftarrow 0$

**function** AllocateRound():
    // Compute remaining budget
    $\varepsilon_{spent} \leftarrow \text{ComputeEpsilon}(\lambda, \delta_{total})$
    $\varepsilon_{remaining} \leftarrow \varepsilon_{total} - \varepsilon_{spent}$
    // Check if budget allows another round
    **if** $\varepsilon_{remaining} < \varepsilon_{min}$ **then**
        **raise** BudgetExhaustedError
    // Allocate per-round budget
    $\varepsilon_t \leftarrow \varepsilon_{remaining}/(T - \text{rounds\_completed})$
    **return** $\varepsilon_t$

**function** RecordRound($\sigma$, $q$):
    // Update moments after each round
    **for** order $= 1$ to MAX_ORDER **do**
        $\lambda[\text{order}] += \text{ComputeMoment}(\text{order}, \sigma, q)$
    rounds_completed $+= 1$



Fig. 8. Per-client training time per round. Larger hospitals (Berlin: 500 samples) exhibit slightly longer training times. The adaptive training engine compensates by adjusting batch sizes for stragglers.



Fig. 9. Client participation matrix (50 rounds × 5 clients). Participation rates: IT 88%, DE 86%, FR 86%, ES 88%, NL 92%. The framework tolerates 10–15% dropout per round while maintaining convergence.

## I. Hospital Data Distribution

Figure 7 illustrates the non-IID nature of data across the five simulated hospitals. Each hospital exhibits distinct demographic characteristics reflecting real-world geographical variation in European patient populations.
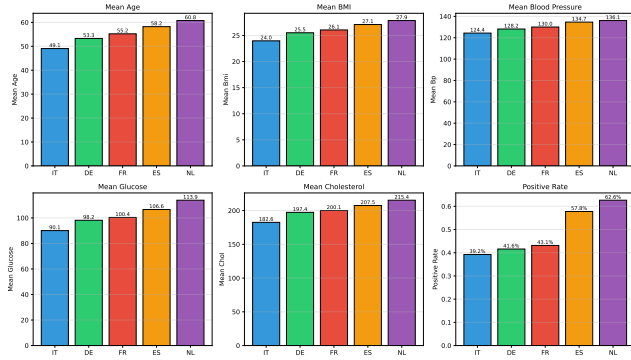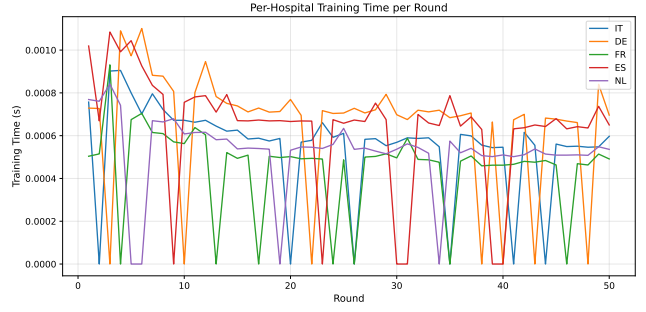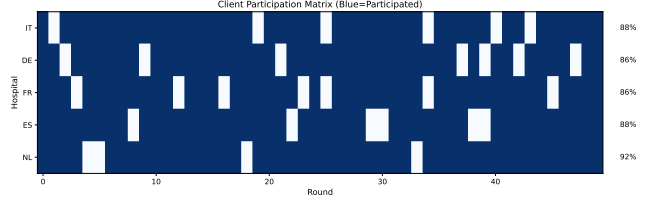


Fig. 7. Data distribution across hospitals. Metrics shown: sample count, mean age, BMI, systolic BP, glucose, and positive class rate. Notable heterogeneity: Amsterdam shows older population (60.8 years mean age) with higher positive rate (62.6%) compared to Rome (49.1 years, 39.2%).

## J. Per-Client Training Time

Figure 8 shows training time variation across clients per round. Differences arise from local dataset sizes (300–500 records), hardware capabilities, and network conditions.

## K. Client Participation Matrix

Figure 9 presents the client participation matrix over 50 training rounds. Not all clients participate in every round due to availability, connectivity, or straggler timeout policies.

## L. Gradient Norm Evolution

Figure 10 tracks gradient L2-norms throughout training. Decreasing gradient norms indicate model convergence; divergent norms suggest instability.
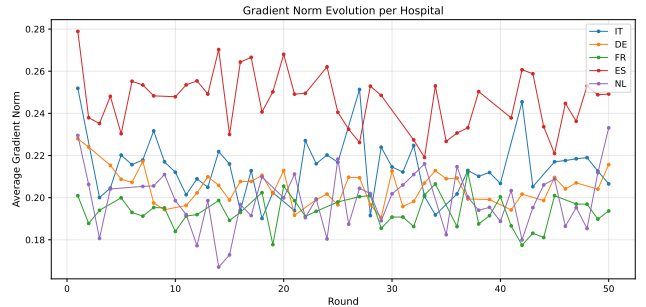


Fig. 10. Gradient norm evolution per client over 50 rounds. All clients show decreasing trends indicating stable convergence. Clipping threshold $C = 1.0$ bounds extreme values for DP compatibility.

## M. Communication Cost Analysis

Figure 11 analyzes per-round communication overhead. For logistic regression with 6 parameters, each gradient transmission is approximately 2.3 KB (32-bit floats + protocol overhead).

## N. Learning Rate Sensitivity

Figure 12 compares convergence across learning rates $\eta \in \{0.01, 0.05, 0.1, 0.2, 0.5\}$. Optimal performance is achieved at $\eta = 0.1$–$0.2$.
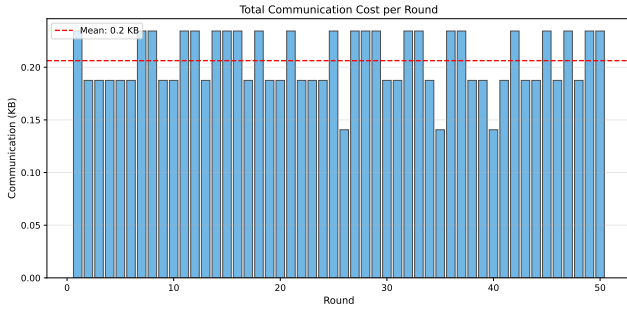
Fig. 11. Cumulative communication cost per round. Linear scaling with participating clients (3.5 KB/client/round). Total 50-round overhead: 875 KB for 5 clients—feasible even for bandwidth-constrained environments.
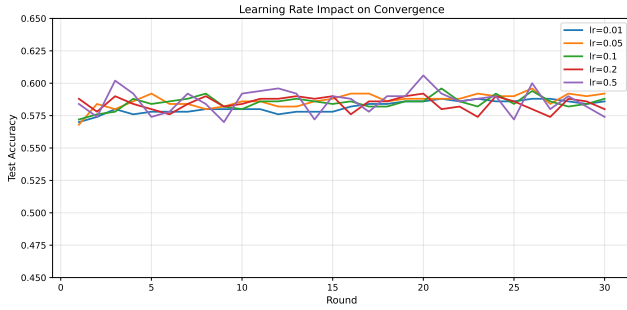


Fig. 14. Per-client accuracy over training rounds. Variance reflects non-IID data: NL (older, higher-risk population) reaches 64% accuracy while FR (mid-range demographics) stabilizes at 55%.



Fig. 12. Learning rate sensitivity analysis. $\eta = 0.01$: slow convergence (53.8% at round 50). $\eta = 0.1$: optimal (58.6%). $\eta = 0.5$: instability with oscillations.

## O. Batch Size Impact

Figure 13 evaluates the effect of batch sizes $\{8, 16, 32, 64, 128\}$ on convergence speed and final accuracy.
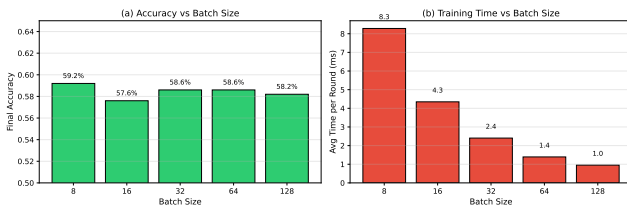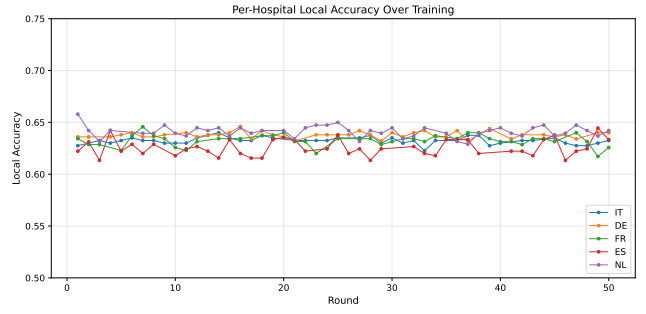


Fig. 13. Batch size impact on convergence. Smaller batches (8–16) provide noisier gradients but faster initial progress. Batch size 32 balances gradient quality and computational efficiency.

## P. Per-Client Accuracy Trajectories

Figure 14 shows individual client accuracy trajectories, revealing heterogeneity in local model performance.

This appendix provides comprehensive benchmarking of five FL algorithms across varying degrees of data heterogeneity. Results inform algorithm selection for different EHDS deployment scenarios.

## Q. Algorithms Evaluated

We compare five foundational FL algorithms, plus seven recent methods (2022–2025) implemented in the framework:

- **FedAvg** [12]: Baseline federated averaging with weighted aggregation.
- **FedProx** [13]: Proximal term ($\mu$) regularizes local drift in non-IID settings.
- **SCAFFOLD** [23]: Variance reduction through control variates correcting client drift.
- **FedAdam** [26]: Adaptive server-side optimization with momentum and second-moment estimates.
- **FedYogi** [26]: Variant of FedAdam with improved stability for sparse gradients.

*Recent algorithms (2022–2025)* extending the framework beyond foundational methods:

- **FedLC** [38]: Calibrates logits via class-frequency margins; composable with any algorithm. Addresses label distribution skew across hospitals.
- **FedSAM** [37]: Sharpness-Aware Minimization for flat minima; drop-in local optimizer replacement.
- **FedDecorr** [39]: Decorrelation regularizer preventing dimensional collapse; composable, negligible overhead.
- **FedSpeed** [40]: Unifies proximal correction with gradient perturbation; fewer communication rounds.
- **FedExP** [41]: Adaptive server step size via POCS analogy; zero client-side changes, hyperparameter-free.
- **FedLESAM** [44]: Extends FedSAM with locally-estimated global perturbation direction ($\theta_g^{(t-1)} - \theta_g^{(t)}$); ICML 2024 Spotlight.
- **HPFL** [45]: Hot-Pluggable FL—aggregates shared backbone only, keeps classifier heads local per client; ICLR 2025.

## R. Non-IID Configuration

Data heterogeneity is controlled via Dirichlet distribution with concentration parameter $\alpha$:

- $\alpha = 0.1$: **Extreme non-IID** — highly skewed label distributions per client
- $\alpha = 0.5$: **High non-IID** — significant heterogeneity
- $\alpha = 1.0$: **Moderate non-IID** — balanced heterogeneity
- $\alpha = 10.0$: **Near-IID** — approximately uniform distributions

## S. Convergence at Different Heterogeneity Levels

Figure 15 shows convergence trajectories for all algorithms across four non-IID levels. Key observations:
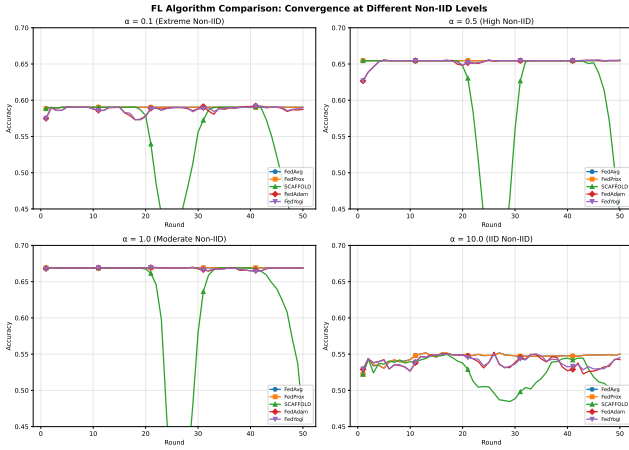


Fig. 15. Algorithm convergence across non-IID levels ($\alpha \in \{0.1, 0.5, 1.0, 10.0\}$). SCAFFOLD and adaptive methods (FedAdam, FedYogi) show superior stability under extreme heterogeneity ($\alpha = 0.1$).

**Findings:**

1) At $\alpha = 0.1$ (extreme non-IID), **SCAFFOLD** achieves most stable convergence due to variance reduction.
2) **FedProx** provides marginal improvement over FedAvg at moderate heterogeneity ($\alpha = 0.5$–$1.0$).
3) **Adaptive methods** (FedAdam, FedYogi) excel in near-IID settings ($\alpha = 10$) but may oscillate under extreme heterogeneity.
4) **FedAvg** remains competitive in near-IID conditions, making it suitable for homogeneous federations.

## T. Final Accuracy vs. Data Heterogeneity

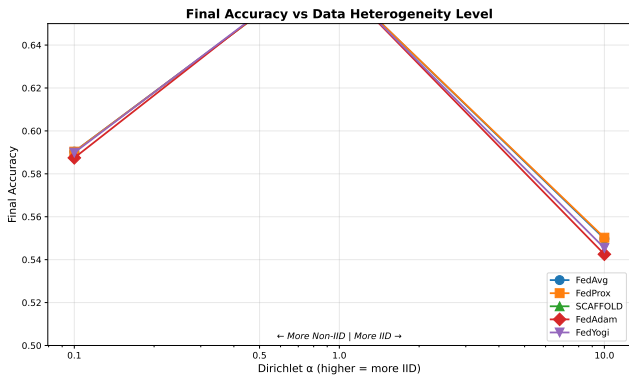Figure 16 summarizes final accuracy (round 50) as a function of heterogeneity level.



Fig. 16. Final accuracy vs. Dirichlet $\alpha$. All algorithms degrade under extreme non-IID conditions. SCAFFOLD shows smallest performance gap between $\alpha = 0.1$ and $\alpha = 10$.

**EHDS Implications:** Cross-border federations with heterogeneous patient populations (different age distributions, disease prevalence, clinical practices) should prefer SCAFFOLD

or FedProx over baseline FedAvg. For deep learning deployments where hospitals exhibit different disease prevalences, FedLC's logit calibration and FedDecorr's decorrelation regularizer provide composable improvements orthogonal to the choice of base aggregation algorithm.

## U. Convergence Speed Analysis

Figure 17 analyzes two convergence metrics: (a) rounds required to reach 55% accuracy threshold, and (b) best accuracy achieved within the first 20 rounds.
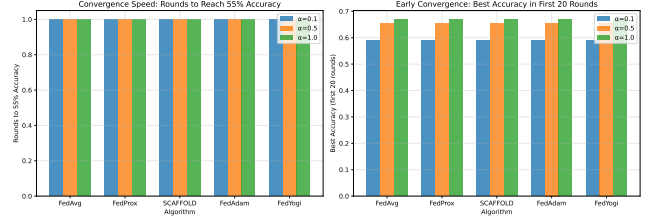


Fig. 17. Convergence speed comparison. Left: rounds to 55% accuracy. Right: best accuracy in first 20 rounds. Adaptive methods converge faster in early rounds but may plateau.

**Resource-Constrained Scenarios:** For EHDS deployments with limited communication budgets or time constraints, adaptive methods (FedAdam, FedYogi) may be preferred for their faster early convergence, despite potential instability in later rounds.

## V. Algorithm Selection Guidelines for EHDS

Table XIII provides recommendations based on federation characteristics:

TABLE XIII
ALGORITHM SELECTION GUIDELINES FOR EHDS DEPLOYMENTS

| Scenario | Recommended | Rationale |
|---|---|---|
| Homogeneous MS | FedAvg | Simplicity, proven |
| Heterogeneous MS | SCAFFOLD | Variance reduction |
| Resource-limited | FedAdam | Fast convergence |
| Privacy-critical | FedAvg + DP | Well-studied DP bounds |
| Sparse participation | FedProx | Dropout resilience |
| Label-imbalanced | FedLC | Class-freq. calibration |
| Deep models, non-IID | FedDecorr | Prevents dim. collapse |
| Comm.-constrained | FedSpeed | Fewer rounds needed |
| No client changes | FedExP | Server-side only |
| SAM + global drift | FedLESAM | Globally-guided flatness |
| Per-hospital classifiers | HPFL | Local decision boundaries |

**MS** = Member States. Heterogeneity arises from demographic differences, clinical coding practices, EHR system variations, and population health profiles across EU regions.

The FL-EHDS framework implements several advanced FL paradigms addressing specific challenges in cross-border

healthcare analytics. This appendix documents these components available in the reference implementation.

## W. Vertical Federated Learning

Vertical FL (VFL) addresses scenarios where different institutions hold different features for the same patients—common in healthcare where hospitals, labs, and pharmacies maintain complementary records.

**Private Set Intersection (PSI):** Before training, participating institutions must identify common patients without revealing their full patient lists. The framework implements RSA-based PSI with the following properties:

- **Privacy**: Neither party learns patients not in the intersection
- **Efficiency**: $O(n \log n)$ complexity for $n$ records
- **EHDS Compliance**: Pseudonymized identifiers only

**Split Learning:** Neural networks are partitioned across institutions. Each party computes forward passes on their features up to a "cut layer," exchanges activations (not raw data), and backpropagates gradients. The framework supports:

- U-shaped splits (features split, labels at one party)
- Configurable cut layer selection
- Gradient compression for bandwidth optimization

---

**Algorithm D.1: Split Learning Forward Pass**
**Input:** Features $X_A$, $X_B$ at parties A, B; cut layer $k$
**Output:** Prediction $\hat{y}$

// Party A: forward to cut layer
$h_A \leftarrow f_{1:k}^A(X_A)$      // Local features → activations
// Party B: forward to cut layer
$h_B \leftarrow f_{1:k}^B(X_B)$      // Local features → activations
// Server: aggregate and complete forward pass
$h \leftarrow \text{Concat}(h_A, h_B)$
$\hat{y} \leftarrow f_{k+1:L}(h)$      // Cut layer → output
**return** $\hat{y}$

---

## X. Byzantine-Resilient Aggregation

In cross-border federations, some participants may be compromised, malfunctioning, or malicious. Byzantine-resilient aggregation protects model integrity against up to $f < n/3$ adversarial clients.

**Implemented Defenses:**

- **Krum**: Selects the gradient closest to its $n-f-2$ nearest neighbors, excluding outliers
- **Trimmed Mean**: Discards the $\beta$ fraction of extreme values per coordinate before averaging
- **Coordinate-wise Median**: Robust estimator immune to single-coordinate attacks
- **Bulyan**: Two-stage defense combining Krum selection with trimmed mean
- **FLTrust**: Server maintains a small trusted dataset to validate client updates

---

**Algorithm D.2: Krum Byzantine Defense**
**Input:** Gradients $\{g_1, \ldots, g_n\}$, Byzantine bound $f$
**Output:** Selected gradient $g^*$

**for each** gradient $g_i$ **do**
     // Compute distances to all other gradients
     $D_i \leftarrow \{\|g_i - g_j\|^2 : j \neq i\}$
     // Sum of $n - f - 2$ smallest distances
     $s_i \leftarrow \sum_{d \in \text{smallest}_{n-f-2}(D_i)} d$
// Select gradient with minimum score
$g^* \leftarrow g_{\arg\min_i s_i}$
**return** $g^*$

---

## Y. Continual Federated Learning

Healthcare data evolves over time—new diseases emerge, treatment protocols change, and patient demographics shift. Continual FL addresses catastrophic forgetting while adapting to distribution drift.

**Implemented Strategies:**

- **Elastic Weight Consolidation (EWC)**: Protects important weights using Fisher Information as importance measure
- **Learning without Forgetting (LwF)**: Knowledge distillation from previous model prevents forgetting
- **Experience Replay**: Maintains memory buffer of representative past samples
- **Drift Detection**: Statistical tests (ADWIN, Page-Hinkley) detect concept drift and trigger adaptation

The EWC loss function adds a quadratic penalty for deviating from important parameters:

$$\mathcal{L}_{EWC}(\theta) = \mathcal{L}(\theta) + \frac{\lambda}{2} \sum_i F_i(\theta_i - \theta_i^*)^2$$

where $F_i$ is the Fisher Information for parameter $i$ and $\theta^*$ are the optimal parameters for previous tasks.

## Z. Multi-Task Federated Learning

Different hospitals may have different prediction objectives (e.g., mortality, readmission, length-of-stay). Multi-Task FL learns shared representations while allowing task-specific outputs.

**Architectures:**

- **Hard Parameter Sharing**: Common feature extractor, task-specific heads
- **Soft Parameter Sharing**: Separate networks with regularization encouraging similarity
- **FedMTL**: Learns task relationships dynamically during training

## . Hierarchical Federated Learning

The FL-EHDS framework implements a four-tier hierarchy reflecting EU governance structure:

1) **Client Tier**: Individual hospitals/data holders
2) **Regional Tier**: Regional aggregators (e.g., Lombardy, Bavaria)

3) **National Tier**: National HDABs coordinate Member State aggregation
4) **EU Tier**: HealthData@EU central aggregator

This hierarchy reduces communication costs (hospitals communicate with regional aggregators, not directly with EU server) and aligns with EHDS governance where HDABs have national jurisdiction.

### . *Personalized Federated Learning*

Standard FL produces a single global model, but healthcare institutions have distinct patient populations requiring personalized predictions. Personalized FL learns models adapted to each client's data distribution while benefiting from collaborative training.

**Implemented Approaches:**

- **FedPer (Federated Personalization)**: Shares base layers globally, keeps personalization layers local. Suitable when institutions have similar low-level features but different high-level patterns.
- **pFedMe (Personalized FedAvg with Moreau Envelopes)**: Uses Moreau envelopes to decouple personalization from global model learning. Provides theoretical convergence guarantees.
- **Per-FedAvg (Personalized FedAvg)**: Meta-learning approach where the global model serves as initialization for local fine-tuning. Based on MAML (Model-Agnostic Meta-Learning).
- **APFL (Adaptive Personalized FL)**: Learns mixing coefficient $\alpha$ between global and local models per client. Automatically balances personalization vs. generalization.
- **Ditto**: Adds personalization regularization term to local objectives, ensuring local models remain close to global while adapting to local data.

---

**Algorithm D.3: pFedMe Local Update**

**Input:** Local data $\mathcal{D}_k$, global model $\theta$, personal model $\theta_k$, $\lambda$, $\eta$

**Output:** Updated personal model $\theta'_k$

// Personalization step: optimize local objective
**for** $i = 1$ to $R$ **do**
    $\theta_k \leftarrow \theta_k - \eta \nabla \mathcal{L}(\theta_k; \mathcal{D}_k)$
// Moreau envelope update: balance with global
$\theta'_k \leftarrow \theta_k - \lambda(\theta_k - \theta)$
// Compute gradient for global model
$g_k \leftarrow \lambda(\theta - \theta'_k)$
**return** $\theta'_k$, $g_k$

---

**EHDS Relevance:** Personalized FL is critical for EHDS because Member States have different healthcare systems, disease prevalence, and clinical practices. A model trained for Finnish cardiovascular risk may not transfer directly to Italian populations due to dietary, genetic, and healthcare access differences.

### . *Asynchronous Federated Learning*

Synchronous FL requires all clients to complete local training before aggregation, creating bottlenecks when clients have heterogeneous computational resources or network connectivity. Asynchronous FL allows clients to contribute updates independently.

**Implemented Strategies:**

- **FedAsync**: Server aggregates each client update immediately upon arrival. Uses staleness-weighted averaging to discount outdated gradients.
- **FedBuff (Buffered Asynchronous)**: Collects $K$ updates in a buffer before aggregating. Balances asynchrony benefits with aggregation quality.
- **ASO-Fed (Adaptive Synchronization)**: Dynamically adjusts synchronization frequency based on client staleness distribution.
- **Semi-Asynchronous**: Waits for a fraction $\alpha$ of clients (e.g., 70%) before aggregating, bounding maximum staleness.

---

**Algorithm D.4: FedAsync with Staleness Weighting**

**Input:** Client update $\Delta_k$, client round $t_k$, server round $t$

**Output:** Updated global model $\theta$

// Compute staleness
$\tau \leftarrow t - t_k$       // How many rounds old is this update?
// Staleness-weighted coefficient (polynomial decay)
$\alpha \leftarrow (1 + \tau)^{-a}$       // $a > 0$ controls decay rate
// Apply weighted update
$\theta \leftarrow \theta + \alpha \cdot \eta \cdot \Delta_k$
**return** $\theta$

---

**Staleness Functions:**

- **Constant**: $\alpha(\tau) = 1$ (ignore staleness)
- **Polynomial**: $\alpha(\tau) = (1 + \tau)^{-a}$
- **Exponential**: $\alpha(\tau) = e^{-a\tau}$
- **Hinge**: $\alpha(\tau) = 1$ if $\tau \leq \tau_{max}$, else 0

**EHDS Relevance:** Cross-border federations span institutions with vastly different IT infrastructure. A Finnish hospital with modern HPC may complete training in minutes, while a rural clinic in another Member State may require hours. Asynchronous FL prevents fast clients from waiting for stragglers.

### . *Fairness-Aware Federated Learning*

FL can amplify existing healthcare disparities if minority populations or under-resourced institutions are underrepresented in training. Fairness-aware FL ensures equitable model performance across demographic groups and participating institutions.

**Fairness Dimensions:**

- **Client Fairness**: Equitable accuracy across participating hospitals
- **Group Fairness**: Balanced performance across demographic groups (age, gender, ethnicity)
- **Individual Fairness**: Similar predictions for similar patients

**Implemented Approaches:**

- **q-FedAvg**: Reweights client contributions to minimize worst-case client loss. Parameter $q$ controls fairness-accuracy tradeoff.

- **AFL (Agnostic Federated Learning)**: Optimizes for the worst-performing distribution mixture, providing robustness guarantees.
- **FedMGDA+ (Multi-Gradient Descent)**: Finds Pareto-optimal update direction balancing all client objectives.
- **TERM (Tilted Empirical Risk Minimization)**: Uses tilted losses to emphasize underperforming groups.
- **FairFed**: Adds fairness constraints (demographic parity, equalized odds) to the federated optimization.

---

**Algorithm D.5: q-FedAvg Fair Aggregation**

**Input:** Client losses $\{L_1, \ldots, L_K\}$, updates $\{\Delta_1, \ldots, \Delta_K\}$, $q$
**Output:** Fair aggregated update $\Delta$

// Compute fairness weights (emphasize high-loss clients)
**for each** client $k$ **do**
    $w_k \leftarrow L_k^q$            // Higher loss $\rightarrow$ higher weight
// Normalize weights
$W \leftarrow \sum_k w_k$
**for each** client $k$ **do**
    $w_k \leftarrow w_k/W$
// Weighted aggregation
$\Delta \leftarrow \sum_k w_k \cdot \Delta_k$
**return** $\Delta$

---

**Fairness Metrics:**
- **Performance Variance**: $\text{Var}(\{L_k\}_{k=1}^K)$ across clients
- **Worst-case Loss**: $\max_k L_k$
- **Demographic Parity Gap**: $|P(\hat{Y} = 1|A = 0) - P(\hat{Y} = 1|A = 1)|$
- **Equalized Odds Gap**: Difference in TPR/FPR across groups

**EHDS Relevance:** The EHDS explicitly aims to benefit all European citizens. Models that perform well on average but poorly for specific populations (e.g., rare diseases, elderly patients, specific ethnic groups) violate the regulation's equity principles. Fairness-aware FL ensures no Member State or demographic group is systematically disadvantaged.

Production deployment of FL-EHDS requires enterprise-grade infrastructure. This appendix documents the infrastructure modules available in the reference implementation.

### . Communication Layer

The framework provides two transport options optimized for different deployment scenarios:

**gRPC (Google Remote Procedure Call):**
- Bidirectional streaming for continuous gradient exchange
- Protocol Buffers serialization (30% bandwidth reduction vs. JSON)
- HTTP/2 multiplexing reduces connection overhead
- Ideal for: data center deployments, low-latency requirements

**WebSocket:**
- Browser-compatible for web-based FL clients
- Firewall-friendly (standard HTTP upgrade)
- Event-driven architecture for real-time updates
- Ideal for: edge deployments, browser-based participation

---

**Communication Manager Configuration**

```
transport: gRPC | WebSocket
compression: gzip | lz4 | zstd | none
chunk_size: 1MB              // For streaming large models
retry_policy:
  max_retries: 3
  backoff: exponential
  base_delay: 1s
connection_pool:
  max_connections: 100
  idle_timeout: 300s
```

---

### . Serialization Layer

Efficient model serialization is critical for cross-border FL where bandwidth is constrained. The framework implements:

**Protocol Buffer-style Binary Format:**
- Tensor metadata (shape, dtype) + raw binary data
- 30% smaller than JSON, 15% smaller than pickle
- Cross-platform compatibility (Python, C++, Java)

**Delta Serialization:**
- Transmits only changed parameters between rounds
- Sparse encoding for gradient updates
- Up to 90% bandwidth reduction for fine-tuning scenarios

**EHDS-Compliant Serialization:**
- Embeds permit ID, timestamp, and provenance metadata
- Cryptographic signatures for integrity verification
- Audit trail fields for GDPR Article 30 compliance

### . Caching Layer

Distributed caching accelerates FL operations and provides fault tolerance:

**Redis-based Implementation:**
- **Model Checkpoints**: Global model snapshots for recovery
- **Client State**: Local model states, optimizer moments
- **Metrics**: Real-time training metrics for monitoring

**Features:**
- LRU/LFU/TTL eviction policies
- Distributed locking for concurrent access
- Automatic serialization/deserialization
- Cache warming for predictable latency

---

**Algorithm E.1: Distributed Lock for Aggregation**

**Input:** Lock name, TTL, client ID
**Output:** Lock acquired (boolean)

// Attempt to acquire lock atomically
acquired $\leftarrow$ Redis.SET(lock_name, client_id, NX, EX=TTL)
**if** acquired **then**
    // Perform aggregation
    PerformAggregation()
    // Release lock only if we own it
    **if** Redis.GET(lock_name) == client_id **then**
        Redis.DEL(lock_name)
**return** acquired

*. Orchestration Layer*

Enterprise FL deployments require automated scaling and management:

**Kubernetes Integration:**

- Deploys FL clients/aggregators as Kubernetes pods
- Horizontal Pod Autoscaler (HPA) for elastic scaling
- ConfigMaps for FL hyperparameter injection
- Secrets for credential management (HDAB API keys)

**Ray Distributed Computing:**

- Actor-based FL client/aggregator implementation
- Automatic resource management and fault tolerance
- Integration with Ray Tune for federated HPO
- Object store for efficient gradient sharing

**Auto-Scaling Policies:**

- **Reactive**: Scale based on queue depth/latency
- **Predictive**: ML-based workload forecasting
- **Scheduled**: Time-based scaling for known patterns

*. Monitoring Layer*

Production observability follows the three pillars: metrics, logs, and traces.

**Prometheus Metrics:**

- **Counters**: rounds_total, permits_validated, gradients_received
- **Gauges**: active_clients, global_model_loss, privacy_budget_remaining
- **Histograms**: round_duration, communication_latency, aggregation_time
- **Summaries**: gradient_norm_quantiles

**Grafana Dashboards:** The framework auto-generates Grafana dashboard JSON for:

- FL training progress (loss, accuracy per round)
- Client participation and health
- Communication latency heatmaps
- Privacy budget consumption
- EHDS compliance status

**Alerting:**

- Privacy budget exhaustion warning
- Client dropout rate exceeding threshold
- Model divergence detection
- Permit expiration alerts

*. Model Watermarking*

Intellectual property protection for FL models trained on EHDS data:

**Watermarking Techniques:**

- **Spread Spectrum**: Embeds watermark in frequency domain, robust to fine-tuning
- **LSB (Least Significant Bit)**: Embeds in low-order bits of weights
- **Backdoor-based**: Specific input-output pairs as ownership proof
- **Passport Layers**: Dedicated layers encode ownership information

**EHDS Use Case:** Watermarking enables tracking of model provenance—which data permits, which hospitals contributed, and usage restrictions. This supports EHDS requirements for accountability in secondary use.

*. Cross-Silo Enhancements*

Enterprise FL deployments require advanced capabilities beyond basic aggregation. The framework implements three cross-silo enhancements for production EHDS deployments.

**Multi-Model Federation:** Ensemble learning across federated models improves robustness and generalization:

- **Weighted Voting**: Performance-based weighting of ensemble members
- **Stacking**: Meta-learner trained on base model outputs
- **Mixture of Experts**: Input-dependent gating selects expert models
- **Diversity Enforcement**: Ensures ensemble members capture different patterns

**Automatic Algorithm Selection:** Data-driven selection of the optimal FL algorithm:

- **Task Analysis**: Automatic detection of data heterogeneity (IID vs. non-IID)
- **Multi-Armed Bandit**: UCB/Thompson Sampling for exploration-exploitation
- **Criterion-based**: Optimize for accuracy, convergence, fairness, or privacy

**Adaptive Aggregation:** Dynamic switching between aggregation algorithms based on runtime metrics:

---

**Algorithm E.2: Adaptive Aggregation**

**Input:** Client updates, metrics history, cooldown period
**Output:** Aggregated model, selected algorithm

// Evaluate current algorithm performance
score ← WeightedScore(loss, accuracy, variance, convergence)
// Check if switch would be beneficial
**if** RoundsSinceSwitch > CooldownPeriod **then**
   **for each** candidate ∈ {FedAvg, FedProx, SCAFFOLD, ...}
**do**
     alt_score ← EstimatePerformance(candidate)
     **if** alt_score > score + SwitchThreshold **then**
       SwitchTo(candidate)
// Perform aggregation with current algorithm
aggregated ← CurrentAlgorithm.Aggregate(updates)
**return** aggregated

---

**Supported Algorithms:**

- **Standard**: FedAvg, FedProx, SCAFFOLD, FedNova
- **Adaptive**: FedAdam, FedYogi, FedAdagrad
- **Byzantine-resilient**: Krum, TrimmedMean, Median, Bulyan

**EHDS Relevance:** Cross-border federations exhibit dynamic characteristics: client availability varies by time zone, data heterogeneity depends on participating Member States, and regulatory requirements may change. Adaptive aggregation automatically adjusts to these conditions without manual intervention.

Beyond the FHIR R4 preprocessing described in Algorithm 7, the FL-EHDS framework implements additional interoperability standards required for cross-border health data exchange.

### *. OMOP Common Data Model*

The Observational Medical Outcomes Partnership (OMOP) CDM provides a standardized analytical format for observational health data. Many European research networks (e.g., EHDEN, OHDSI) use OMOP.

**Framework Support:**

- **ETL Pipelines**: Transform source EHR to OMOP CDM v5.4
- **Vocabulary Mapping**: Standard concepts (SNOMED, ICD10, LOINC, RxNorm)
- **Cohort Definitions**: ATLAS-compatible cohort SQL generation
- **Feature Extraction**: FeatureExtraction package integration for ML-ready datasets

**FL Integration:** OMOP standardization enables consistent feature engineering across sites:

1) Each hospital transforms local EHR to OMOP
2) Feature extraction produces identical schema across sites
3) FL training proceeds on homogeneous feature spaces

### *. IHE Integration Profiles*

Integrating the Healthcare Enterprise (IHE) profiles ensure secure, auditable data exchange:

**Implemented Profiles:**

- **ATNA (Audit Trail and Node Authentication)**:
  - TLS mutual authentication between FL nodes
  - Syslog audit messages for all data access events
  - RFC 5424 compliant audit record format
- **BPPC (Basic Patient Privacy Consents)**:
  - Maps Article 71 opt-out to BPPC consent documents
  - XDS.b integration for consent document sharing
  - Consent enforcement at FL training initiation
- **XCA (Cross-Community Access)**:
  - Cross-border document query and retrieve
  - Initiating/Responding Gateway implementation
  - Patient identity correlation across communities
- **PIX/PDQ (Patient Identifier Cross-referencing/Demographics Query)**:
  - Patient matching across institutional boundaries
  - Pseudonymization-aware identity management
  - Integration with national eHealth infrastructures
- **XUA (Cross-Enterprise User Assertion)**:
  - SAML 2.0 assertions for federated authentication
  - Role-based access control integration
  - HDAB authorization token propagation

### *. Cross-Border Data Exchange*

The framework implements the technical specifications for HealthData@EU cross-border exchange:

**Message Formats:**

- EHDS Data Permit Exchange Format (JSON-LD)
- Federated Query Protocol (based on SPARQL Federation)
- Model Update Message Format (Protocol Buffers)

**Security Requirements:**

- eIDAS-compliant electronic signatures for permits
- TLS 1.3 for all cross-border communication
- Certificate-based node authentication (EU trust framework)

**Metadata Standards:**

- DCAT-AP Health extension for dataset cataloging
- Provenance metadata (W3C PROV-O)
- Data quality indicators per EMA guidelines

### *. Reference Implementation Architecture*

Figure 18 illustrates how the interoperability components integrate within the FL-EHDS framework.
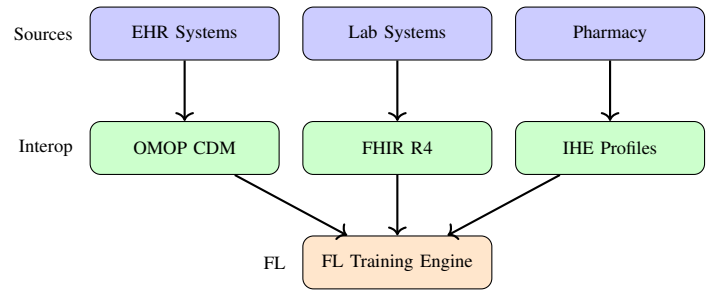


Fig. 18. Interoperability layer integrating heterogeneous data sources for FL training.

This appendix provides extended details for the clinical imaging experiments reported in Section V. Main results are in Table XI.

### *. Datasets*

Three clinical imaging datasets cover representative EHDS secondary use scenarios:

- **Chest X-ray** [31]: 5,860 pediatric radiographs (NORMAL/PNEUMONIA, 2.7:1 imbalance). Binary classification for pneumonia detection.
- **Brain Tumor MRI**: 3,064 T1-weighted CE MRI slices. 3-class classification (glioma, meningioma, pituitary tumor).
- **Skin Cancer**: 3,297 dermoscopy images. Binary classification (benign vs. malignant).

### *. Model Architectures*

Two architectures are available in the framework:

- **HealthcareResNet**: ResNet-18 [43] pretrained on ImageNet, with GroupNorm replacing BatchNorm for FL stability. Partial backbone freezing (level 1: conv1+bn1

frozen). FedBN [33] skips normalization layers during aggregation. $\sim$11.2M parameters.

- **HealthcareCNN**: 5-block CNN with Group-Norm, progressive channels $(32{\rightarrow}512)$, graduated Dropout $(0.15{\rightarrow}0.3)$. Classifier: Flatten$\rightarrow$FC(512)$\rightarrow$FC(128)$\rightarrow$FC($K$). $\sim$12M parameters.

Data augmentation: random horizontal flip, rotation ($\pm15°$), brightness jitter ($\pm10\%$). ImageNet normalization.

### . *Experimental Configuration*

The V2 imaging configuration uses improved hyperparameters based on preliminary ablation:

- 5 hospitals, 25 rounds, 3 local epochs, batch size 32
- Adam optimizer (lr=0.001), early stopping (patience=6)
- Non-IID via Dirichlet $\alpha$=0.5
- FedBN enabled, partial backbone freeze (level 1)
- 7 algorithms: FedAvg, FedProx, Ditto, FedLC, FedExP, FedLESAM, HPFL
- 3 seeds per configuration (42, 123, 456)

### . *Reproducibility*

All experiments are reproducible via the reference implementation:

```
cd fl-ehds-framework
# Full V2 experiments (7 algo x 5 datasets x 3 seeds)
python -m benchmarks.run_full_experiments
# Quick validation (˜1-2h)
python -m benchmarks.run_full_experiments --quick
# Resume after interruption
python -m benchmarks.run_full_experiments --resume
```

Results, checkpoints, and logs are auto-saved to benchmarks/paper_results/.