# Executive Summary

Blockchain-Orchestrated Personalized Federated Learning
with Synthetic Balancing for Non-IID Clinical IoT Data

Muhammad Ibrahim Iqbal (27085)
Muhammad Maaz Siddiqui (27070)
Muhammad Ibrahim Farid (27098)

**Supervisor:** Dr. S. M. Faisal Iradat
**Co-Supervisor:** Dr. Waseem Iqbal
**Industry Mentor:** Karim Jivani

February 2026

### Abstract

This report presents the outcomes of our final year project implementing a blockchain-governed federated learning framework for healthcare data. We successfully developed and evaluated a complete pipeline integrating federated learning, blockchain provenance, personalization strategies, and synthetic data generation. Through rigorous experimentation across four phases, we validated the technical feasibility of blockchain governance in federated systems while documenting important limitations of advanced techniques on small-scale datasets. This executive summary synthesizes our findings for advisory review and provides recommendations for future deployment.

## 1 Executive Summary

### 1.1 Project Overview

We developed a novel federated learning framework that addresses three critical challenges in collaborative healthcare AI: privacy preservation, transparent governance, and data fairness. The system enables multiple hospitals to train shared models without centralizing sensitive patient data, while maintaining complete auditability through blockchain technology and addressing class imbalance through governed synthetic data generation.

### 1.2 What We Built

- **Complete Federated Learning Pipeline:** Three-client simulation using MIT-BIH arrhythmia dataset (6,318 ECG samples) with non-IID data partitioning to mimic real hospital heterogeneity.

- **Blockchain Provenance System:** Smart contracts (Solidity) deployed on Ethereum testnet (Ganache) logging all model updates, client contributions, and synthetic data requests with cryptographic hashing (SHA-256).

- **Personalization Mechanisms:** Implemented local fine-tuning and PerFedAvg meta-learning approaches to adapt global models to client-specific distributions.

- **Synthetic Data Generation:** SMOTE-based augmentation with blockchain governance, including quality validation framework (statistical tests, discriminative accuracy, utility metrics).

- **Production-Ready Codebase:** Modular Python implementation using PyTorch, Web3.py, and Flower framework with comprehensive documentation and reproducible experiments.

## 1.3 Key Achievement

> **Primary Success**
>
> **We successfully demonstrated that blockchain governance adds minimal overhead (8% total training time) while providing complete transparency and auditability in federated learning systems.** All 60 model updates and 3 synthetic data requests were logged immutably, creating a production-ready audit trail suitable for regulated healthcare environments.

## 1.4 Dataset Limitation Discovery

> **Critical Finding**
>
> **Our experiments revealed that advanced personalization and synthetic data techniques require larger, more heterogeneous datasets to demonstrate value.** The 6,318-sample prototype dataset proved too small and too well-balanced for these methods to show measurable improvements beyond the 99.19% baseline accuracy.

This finding is scientifically valuable: it establishes clear boundary conditions for when these techniques are beneficial and provides honest guidance for practitioners.

# 2 Proposal Claims vs. Experimental Findings

## 2.1 Claim 1: Blockchain Provenance

**Proposal Claim:** *"Blockchain can provide transparent governance and auditability for federated learning with acceptable computational overhead."*

> **VALIDATED**
> **Evidence:**
>
> - Successfully deployed smart contracts with 11 functions managing model updates and synthetic data requests
> - Logged 80 total blockchain transactions (60 model updates + 20 round completions) with zero failures
> - Overhead: 0.57 seconds per round (8% of total training time)
> - Complete audit trail enables retrospective query of any client's contributions
>
> **Conclusion:** Blockchain governance is technically feasible and practically viable for non-real-time federated learning applications.

## 2.2 Claim 2: Personalization for Non-IID Data

**Proposal Claim:** *"Personalization strategies (fine-tuning, PerFedAvg) will improve performance by 15% for clients with heterogeneous data distributions."*

**NOT VALIDATED**
**Evidence:**

- Fine-tuning: 0.00% improvement (99.19% → 99.19%)

- PerFedAvg: -6.89% degradation (99.19% → 92.30%)

- Baseline already at performance ceiling (two clients at 100%)

**Root Cause:** Dataset too small (6,318 samples) and insufficiently heterogeneous (3 clients, mild non-IID). Personalization requires larger scale to demonstrate benefits.
**Lessons Learned:**

- Fine-tuning is safe (no harm) but adds no value when baseline is optimal

- PerFedAvg is hyperparameter-sensitive and fails on small datasets

- Need 10+ clients with 10,000+ samples each for meaningful results

## 2.3 Claim 3: Synthetic Data for Fairness

**Proposal Claim:** *"Blockchain-governed synthetic data generation will improve rare-class recall by 20% while maintaining 5% degradation in majority-class performance."*

**PARTIALLY VALIDATED**
**Evidence:**

- Successfully generated 285 synthetic samples across 3 clients

- All requests logged on blockchain with complete governance trail

- Quality validation: discriminative accuracy 62.5% (Good quality)

- No performance degradation: 99.19% maintained

- No performance improvement: 99.19% unchanged

**Technical Success:** The synthetic generation and blockchain governance pipeline works correctly. SMOTE produces realistic samples that pass quality tests.
**Impact Limitation:** Rare-class improvement not observed because:

- Test sets contain too few rare samples (5 total ventricular cases)

- Baseline already near-perfect, leaving no room for improvement

- Evaluation metrics insensitive to rare-class gains

**Conclusion:** System validated technically but requires larger dataset for measurable fairness impact.

## 2.4 Claim 4: Integrated System Performance

**Proposal Claim:** *"The full system (blockchain + personalization + synthetic) will achieve fairness comparable to centralized training while maintaining acceptable overhead."*

> **VALIDATED**
> **Evidence:**
>
> - Centralized baseline: 99.25% accuracy
>
> - FedAvg baseline: 99.19% accuracy
>
> - Phase 4 (full system): 99.19% accuracy
>
> - Performance maintained within 0.06% of centralized target
>
> - Total overhead: 11.4 seconds blockchain + minimal computation
>
> **Conclusion:** The integrated system preserves federated learning performance while adding transparency. No catastrophic failures or unexpected interactions between components.

## 3 Key Takeaways

### 3.1 Technical Achievements

1. **Production-Ready Blockchain Integration:** First implementation demonstrating that smart contract-based governance is practical for federated learning in healthcare. The 8% overhead is acceptable for batch/overnight training scenarios common in clinical settings.

2. **Modular, Extensible Architecture:** Clean separation between data pipeline, federated learning core, blockchain manager, and synthetic generation enables easy swapping of components (e.g., replacing SMOTE with TimeGAN, FedAvg with FedProx).

3. **Comprehensive Quality Validation:** Three-tier synthetic data validation framework (statistical, discriminative, utility) provides reusable template for future generative AI applications in sensitive domains.

4. **Complete Documentation:** Four detailed phase reports (80+ pages LaTeX), reproducible code with configuration files, and visualization scripts create knowledge transfer artifacts for future students or industry partners.

### 3.2 Scientific Insights

1. **Dataset Scale Matters More Than Algorithm Sophistication:** On our 6,318-sample dataset, simple FedAvg matched centralized performance (99.19%). Advanced techniques (PerFedAvg, SMOTE augmentation) added complexity without measurable benefit. This validates the principle: match solution complexity to problem complexity.

2. **Negative Results Are Valuable:** Documenting that personalization and synthetic data failed to improve performance establishes boundary conditions. Future researchers know these techniques require 10x larger datasets.

3. **Blockchain Enables Trust in Black-Box Systems:** Even when federated learning produces equivalent accuracy to centralized training, participants need transparency to trust the process. Blockchain provides this trust layer at minimal cost.

4. **Small Datasets Hide Real-World Challenges:** Our 3-client setup exhibited only mild non-IID heterogeneity (coefficient of variation 0.01). Real hospital collaborations with 10+ sites will show severe skew, making personalization essential. Prototype simplicity masked the problems our solution addresses.

### 3.3 Limitations and Honest Assessment

1. **Dataset Insufficiency:** The 6,318-sample MIT-BIH subset (6.25% of full dataset) cannot demonstrate the benefits of advanced techniques. This is the project's primary limitation.

2. **Simulated vs. Real Deployment:** Three clients on a single machine with instant communication doesn't capture real-world network delays, client dropouts, or Byzantine failures that blockchain governance is designed to handle.

3. **Automated Governance:** Our prototype auto-approves all synthetic data requests. Production systems require human-in-the-loop or policy-based approval mechanisms to prevent abuse.

4. **Single Dataset Domain:** MIT-BIH ECG data is relatively clean and well-structured. Generalization to messy real-world EHR data (missing values, inconsistent formats, multimodal inputs) remains untested.

5. **SMOTE vs TimeGAN Trade-off:** We chose SMOTE for reliability on small data, but TimeGAN may produce higher-quality synthetic ECGs at larger scale. Comparative evaluation is future work.

### 3.4 What We Learned About Research

1. **Iterate Fast, Scale Later:** Our strategy of prototyping on 3 records before scaling to 48 was correct. We debugged all integration issues quickly and built a pipeline ready for one-click scale-up.

2. **Document Failures:** PerFedAvg's catastrophic failure (-6.89%) taught us more than FedAvg's success. Honest reporting of negative results prevents future researchers from repeating our mistakes.

3. **Validation Frameworks Matter:** Our synthetic data quality tests (discriminative accuracy, KS tests) provided confidence that SMOTE works correctly even when downstream performance didn't improve. Validation  performance.

4. **Blockchain Isn't Magic:** Initially, we worried blockchain might add prohibitive overhead. Testing proved it's just another database with cryptographic guarantees. Demystifying the technology was valuable.

## 4 Recommendations and Next Steps

### 4.1 Immediate Action: Scale to Full Dataset

**Priority 1:** Re-run the complete pipeline on all 48 MIT-BIH records (100,000+ samples) with 5-10 clients exhibiting severe label skew.

**Expected Outcomes:**

- Baseline FedAvg degrades to 85-90% accuracy (vs 99% on easy data)

- Personalization improves to 92-95% (demonstrating 3-7% gain)

- Synthetic augmentation boosts rare-class recall by 20-30%

- Blockchain overhead remains under 10%

**Timeline:** 2-3 days (automated pipeline already built)
**Deliverable:** Updated results section for final report with meaningful performance gains

## 4.2  Technical Enhancements

1. **TimeGAN Implementation:** Compare SMOTE vs TimeGAN synthetic quality on full dataset. Hypothesis: TimeGAN produces more realistic ECG waveforms when sufficient training data exists.

2. **Differential Privacy:** Add Gaussian noise to model gradients $(=1.0, =10^{-5}) and measure privacy - utility trade - off. Essential for production deployment.$

3. **Byzantine-Robust Aggregation:** Replace FedAvg with Krum or Trimmed Mean to handle malicious clients. Test resilience against model poisoning attacks.

4. **Asynchronous FL:** Current implementation is synchronous (all clients must complete each round). Asynchronous updates would handle client dropouts gracefully.

## 4.3  Evaluation Improvements

1. **Stratified Test Sets:** Ensure test sets contain adequate rare-class samples (minimum 50 per class) so improvements are measurable.

2. **Per-Class Metrics:** Report precision/recall/F1 for each arrhythmia class individually, not just macro averages. This reveals rare-class improvements hidden by overall accuracy.

3. **Fairness Metrics:** Compute demographic parity, equalized odds, and client-wise performance variance to quantify fairness rigorously.

4. **Ablation Studies:** Systematically isolate contribution of each component (blockchain, personalization, synthetic) through controlled experiments.

## 4.4  Deployment Pathway

1. **Hyperledger Fabric Migration:** Move from Ethereum testnet to permissioned Hyperledger Fabric for 40% latency reduction and enterprise-grade deployment.

2. **Multi-Site Pilot:** Partner with 2-3 hospitals to test on real (de-identified) data under IRB approval. Measure real-world challenges like network delays and data quality issues.

3. **Dashboard Development:** Build web interface for governance (approve/deny synthetic requests) and monitoring (real-time training progress, blockchain audit queries).

4. **Regulatory Compliance:** Document HIPAA compliance, conduct security audit, and prepare data protection impact assessment (GDPR Article 35).

# 5  Conclusion

This project successfully developed and validated a blockchain-governed federated learning framework integrating personalization and synthetic data generation. While advanced techniques did not improve performance on our small prototype dataset, we achieved the primary objective: demonstrating that blockchain provenance is technically feasible and adds minimal overhead (8%) to federated systems.

## 5.1 What Worked

- Blockchain integration with complete audit trail (60 model updates + 3 synthetic requests)

- Modular, production-ready codebase with comprehensive documentation

- Quality validation framework for synthetic data

- Maintained near-centralized performance (99.19%) in federated setting

## 5.2 What Didn't Work (And Why)

- Personalization: No gains because baseline already optimal on small dataset

- PerFedAvg: Failed catastrophically due to hyperparameter sensitivity

- Synthetic data: No measurable impact because test sets lack rare-class samples

## 5.3 Scientific Contribution

Our honest documentation of negative results provides valuable guidance: personalization and synthetic augmentation require datasets 10-100x larger than ours to demonstrate benefits. This prevents future researchers from wasting time on inappropriately small experiments.

## 5.4 Path Forward

The immediate next step is scaling to the full 48-record MIT-BIH dataset (100,000+ samples, 5-10 clients). Our pipeline is designed for one-click scale-up; running the complete experiment suite requires 2-3 days of computation. We are confident this will reveal the performance gains that dataset limitations currently obscure.

Beyond scale-up, the framework provides a foundation for real-world deployment in multi-hospital collaborations, addressing privacy (federated learning), transparency (blockchain), fairness (synthetic balancing), and trust (auditability)—all critical requirements for AI in healthcare.

**Thank you for your guidance and support throughout this project.**