

# Privacy-Preserving Vertical Split Learning for Healthcare AI with Prototype-Based Anomaly Detection Against Data Poisoning

Dan Xue

Department of Computer Science, University of Alabama at Birmingham, Birmingham, AL, USA.

xuedan@uab.edu

Niklas Barnett

Department of Computer Science and Engineering, University at Buffalo, Buffalo, NY, USA.

barnett1973@buffalo.edu

## Abstract

The integration of artificial intelligence into healthcare systems promises transformative improvements in diagnostics, treatment personalization, and operational efficiency, yet it simultaneously amplifies concerns regarding patient data privacy and the integrity of learning pipelines. Vertical split learning has emerged as a compelling paradigm that enables multiple healthcare institutions to collaboratively train deep neural networks without sharing raw feature-level data, thereby preserving confidentiality while leveraging complementary data modalities. However, the distributed nature of vertical split learning introduces new attack surfaces, particularly data poisoning and backdoor attacks that can compromise model behavior without violating privacy boundaries. This paper presents a comprehensive system-level analysis of a privacy-preserving vertical split learning framework augmented with prototype-based anomaly detection to counter data poisoning. We examine architectural trade-offs between privacy guarantees, communication overhead, and model accuracy, and we discuss how prototype-based mechanisms, which rely on learning compact class representations, can detect anomalous updates or poisoned samples at the cut layer during split learning. The proposed framework integrates differential privacy mechanisms, secure aggregation protocols, and a prototype consistency verification module that identifies deviations from expected latent distributions. Beyond technical design, we explore governance implications, deployment challenges in heterogeneous hospital networks, regulatory compliance with HIPAA and GDPR, and sustainability considerations regarding computational and energy costs. System-level robustness is evaluated through cross-domain comparisons with horizontal federated learning and fully centralized approaches, highlighting the nuanced benefits and limitations of vertical split learning in sensitive healthcare environments. We conclude with forward-looking perspectives on adaptive defense strategies, standardization of privacy-preserving benchmarks, and the role of policy in fostering trustworthy AI infrastructure.

## Keywords

vertical split learning, privacy-preserving machine learning, healthcare AI, anomaly detection, data poisoning, prototype learning, federated learning, security, governance.

## 1. Introduction

Healthcare institutions increasingly rely on machine learning models to analyze patient data for diagnosis, prognosis, and treatment planning. However, the sensitive nature of medical records imposes stringent privacy requirements that often preclude the centralization of data across hospitals or research networks [9]. Federated learning and its variant, split learning, have been proposed to enable collaborative model training without exposing raw data [1, 11]. Unlike horizontal federated learning where each party holds distinct samples with the same feature space, vertical split learning addresses scenarios where different parties hold different features for the same set of patients, a common situation in healthcare where one hospital may possess imaging data while another holds genomic or laboratory results [2, 3]. In vertical split learning, the neural network is partitioned across two or more parties; each party computes intermediate representations on their local features and sends these activations to a third party (often called the coordinator or active party) that completes the forward propagation and backpropagation. This architecture significantly reduces the amount of information exchanged compared to raw data sharing, but it also creates new vulnerabilities. Adversaries controlling one party can inject poisoned samples or manipulate gradients to implant backdoors, thereby subverting the global model without ever accessing another party's raw data [5, 6]. The threat is particularly acute in healthcare because a compromised model could systematically misdiagnose a condition or fail to detect a critical anomaly, leading to severe clinical consequences. Therefore, robust anomaly detection mechanisms that operate within the split learning protocol are essential. Prototype-based methods, which learn a compact representation of each class as a centroid in the latent space, offer a promising defense by comparing incoming activations or gradient updates against these prototypes to flag outliers [8]. The recent work of Shui et al. introduced ProtoGuard-SL, a prototype consistency based backdoor defense for vertical split learning, demonstrating that such approaches can maintain high detection rates with minimal impact on model performance [15]. In this paper, we adopt a system-level perspective to analyze how prototype-based anomaly detection can be integrated into a privacy-preserving vertical split learning framework for healthcare AI. We examine the structural trade-offs, architectural choices, governance challenges, and deployment realities that must be addressed to make such a system robust, fair, and sustainable.

## 2. Background and Related Work

Vertical split learning draws on the broader taxonomy of distributed machine learning, where the model is partitioned across data silos. Early works by Vepakomma et al. and Gupta and Raskar formalized the split learning paradigm for health applications, demonstrating that it is possible to train deep networks without sharing raw data by exchanging only intermediate activations and gradients [2, 3]. This approach naturally complements privacy-enhancing techniques such as differential privacy, which adds calibrated noise to the shared information to prevent inference of sensitive attributes [4, 13]. However, the interaction between split learning and differential privacy introduces a utility-privacy trade-off: stronger privacy guarantees require larger noise, which degrades model accuracy. Recent studies have explored adaptive noise schedules and local differential privacy to mitigate this degradation while still protecting against membership inference and attribute disclosure [13, 14]. In parallel, data poisoning attacks have been extensively studied in centralized and federated settings. Biggio et al. demonstrated that even small perturbations in training data can cause support vector machines to misclassify specific inputs [5]. Gu et al. later introduced BadNets, showing that backdoors can be embedded in deep learning models by poisoning a small subset of training samples [6]. In the context of vertical split learning, attacks can be launched

by a malicious party that controls either the feature provider or the coordinator. The adversary may tamper with the local data or invert gradients to inject a trigger pattern that the global model learns to associate with a target label. Defenses against such attacks have focused on robust aggregation, anomaly detection in gradient space, and validation using held-out clean data [17]. Prototype-based anomaly detection stands out because it does not require a large clean dataset nor does it assume a specific attack model. Instead, it leverages the geometric structure of the learned latent space: during training, the system maintains prototypical representations for each class, and any update that significantly deviates from the nearest prototype is flagged as suspicious. This approach is computationally efficient and scales well to multiple parties [8, 15].

### **3. System Architecture and Design Considerations**

A privacy-preserving vertical split learning system for healthcare comprises multiple entities: data owners (e.g., hospitals), a coordinator server, and potentially an auditor for compliance. In the typical architecture, each data owner holds a subset of features for a shared set of patient identifiers. The neural network is partitioned vertically: each owner computes a few initial layers on its local features to produce a latent representation, which is then sent to the coordinator. The coordinator concatenates these representations and passes them through the remaining layers to generate predictions and compute the loss. During backpropagation, the coordinator sends gradients back to each owner, who then update their local network parameters. To enforce privacy, the exchanged activations and gradients are perturbed using differential privacy mechanisms, such as adding Gaussian noise calibrated to a privacy budget epsilon [4]. Additionally, secure multi-party computation techniques can be employed to ensure that no single party learns the exact representation of another party's data [12]. The system architecture must also accommodate prototype-based anomaly detection. This requires the coordinator to maintain a dynamic set of prototypical vectors for each class, updated periodically based on the latent representations received from all parties. When a new batch of activations arrives, the coordinator computes the distance from each activation to the nearest prototype. If the distance exceeds a threshold, the sample is considered anomalous and either excluded from training or sent for manual review. The integration of anomaly detection introduces overhead in terms of memory for storing prototypes and computation for distance calculations, but these costs are typically modest compared to the training itself. A critical design decision is whether the anomaly detection module runs on the coordinator or is distributed across parties. Running it centrally simplifies coordination but creates a single point of trust; running it in a decentralized manner requires secure aggregation of prototype distances without revealing individual sample information. The choice depends on the threat model and the trust assumptions among institutions.

### **4. Privacy-Preserving Mechanisms in Vertical Split Learning**

Privacy in vertical split learning is achieved through a combination of data partitioning, differential privacy, and secure computation. Data partitioning inherently limits exposure because raw features never leave each institution's domain; only abstract representations are shared. However, research has shown that intermediate activations can still leak sensitive information, such as patient demographics or even reconstruct original inputs, especially when the attacker has auxiliary knowledge [14]. To address this, differential privacy is applied to the shared activations and gradients. The standard approach involves clipping the norm of each activation vector and adding Gaussian noise with variance proportional to the sensitivity of the function and the desired privacy budget [13]. This noise must be carefully

calibrated: too little noise leaves the system vulnerable to reconstruction attacks, while too much noise degrades model accuracy. In healthcare applications, where diagnostic accuracy is paramount, the privacy budget must be negotiated among stakeholders and potentially enforced by regulatory frameworks like HIPAA or GDPR. Another privacy-enhancing technique is the use of secure aggregation, where the coordinator receives only aggregated contributions from multiple parties, preventing it from isolating individual activations. Secure multi-party computation protocols, such as those based on secret sharing, can further protect against collusion attacks [12]. However, these protocols introduce significant communication and computational overhead, which can be prohibitive in a healthcare setting where hospitals may have limited network bandwidth or computational resources. A trade-off therefore emerges between the strength of privacy guarantees and the feasibility of real-time training. In vertical split learning, the coordinator already sees aggregated activations from multiple parties, but secure aggregation can prevent the coordinator from linking activations back to a specific party. For prototype-based anomaly detection, secure aggregation poses a challenge because the coordinator needs per-sample distances to detect anomalies. One solution is to have each party compute its own prototype distances locally and send only an encrypted flag indicating whether a sample is anomalous, thereby preserving the privacy of the actual distances. The coordinator then aggregates these flags to decide whether to exclude the sample. This approach adds computational overhead but maintains privacy.

## **5. Prototype-Based Anomaly Detection for Data Poisoning Defense**

Prototype-based anomaly detection leverages the intuition that legitimate samples from the same class cluster around a central prototype in the latent space, whereas poisoned samples or those with backdoor triggers deviate from this cluster. The framework proposed by Shui et al., known as ProtoGuard-SL, specifically addresses vertical split learning by maintaining prototypes at the cut layer where activations from different parties are concatenated [15]. During training, the coordinator updates prototypes using an exponential moving average of the latent representations for correctly classified samples. For each incoming sample, the distance to the nearest prototype is computed; if the distance exceeds a dynamic threshold (e.g., the 95th percentile of distances seen so far), the sample is flagged. The flagged sample can then be excluded from gradient updates or used to trigger a manual audit. This approach is particularly effective against backdoor attacks because the trigger pattern, which is usually small and subtle, pushes the representation away from the natural class cluster. The defense does not require knowledge of the specific trigger shape or poisoning strategy, making it robust to a wide range of attacks. Moreover, prototype-based methods are naturally interpretable: clinicians can inspect the most anomalous samples and understand why they were flagged. The system must also handle the case where multiple parties collude. If a malicious party controls several feature providers, they may collectively push the prototype of a target class toward an anomalous region. To mitigate this, the prototypes should be computed using a robust estimator, such as the geometric median instead of the mean, which is less sensitive to outliers. Additionally, the threshold for anomaly detection should be adaptive, increasing over time as the model converges, to avoid false positives from normal but initially distant samples. The integration of this defense into an end-to-end system requires careful engineering. The anomaly detection module must operate in real time during training without introducing significant latency. Since healthcare training datasets are often moderate in size, the overhead is acceptable. However, in scenarios with streaming data (e.g., real-time patient monitoring), the system must be optimized to handle high throughput. The

selection of the prototype update frequency and the distance metric (e.g., Euclidean versus cosine) also affects performance and must be tuned for the specific data modality.

## **6. Structural Trade-offs and System-Level Implications**

The adoption of vertical split learning with prototype-based anomaly detection involves several structural trade-offs that system architects must navigate. First, there is a trade-off between privacy and utility. Stronger differential privacy noise reduces the accuracy of the classifier, which may be unacceptable in high-stakes clinical decisions. Prototype-based defense adds a layer of robustness against poisoning, but the identification and removal of suspicious samples also reduce the effective training set size, potentially harming generalization. However, the net effect is often positive because the removed samples are those most likely to degrade performance. Second, there is a trade-off between communication efficiency and security. Vertical split learning already reduces the amount of data exchanged compared to raw feature sharing, but adding secure aggregation or multi-party computation increases communication rounds and message sizes. In bandwidth-constrained environments (e.g., rural hospitals), this could delay training cycles. The prototype-based detection module, which requires the coordinator to send back information about anomalies, also adds to the communication burden. A possible optimization is to compress the activation vectors before transmission, but compression may interfere with anomaly detection if the latent space geometry is distorted. Third, the system must balance the centralization of the anomaly detection logic against the desire for decentralization. Centralizing detection at the coordinator simplifies auditing and allows global prototype updates, but it creates a central point of failure and trust. Decentralizing detection requires each party to maintain its own prototypes and thresholds, which may lead to inconsistent flagging and reduced detection rates if the parties have non-overlapping feature sets. Hybrid approaches, where the coordinator maintains global prototypes but parties provide local distance summaries, offer a middle ground. Fourth, fairness considerations arise because the prototype-based defense may disproportionately affect minority groups. If a particular demographic group has naturally higher variance in their latent representations (due to underlying biological variability or smaller sample size), they may be more likely to be flagged as anomalous, leading to systematic exclusion from training and subsequent biased predictions. To ensure fairness, the system should incorporate demographic parity checks into the anomaly detection threshold, adjusting for subgroup-specific statistics. These trade-offs are not merely technical; they have direct implications for the deployment of AI in healthcare systems that serve diverse populations.

## **7. Deployment, Governance, and Policy Considerations**

Deploying a privacy-preserving vertical split learning system with prototype-based anomaly detection in real healthcare networks requires careful attention to governance frameworks. Each participating institution must sign data use agreements that specify the privacy budget, the acceptable level of accuracy degradation, and the procedures for handling flagged anomalies. Regulatory compliance with HIPAA in the United States and GDPR in Europe mandates that patient data cannot be re-identified and that individuals have the right to know how their data is used. Vertical split learning, with its limited data sharing, aligns well with these regulations, but the use of differential privacy and prototype-based defense must be documented and auditable. The governance structure should include an independent ethics committee or data protection officer who oversees the training process and reviews anomaly cases. Furthermore, the system must be designed to allow for model updates and retraining

without violating prior privacy constraints. In practice, hospitals may join or leave the federation over time, requiring mechanisms for secure model transfer and forgetting. Prototype-based anomaly detection can assist in this process by identifying unstable representations when new parties join. The computational sustainability of the system is another governance concern. Training large neural networks across multiple sites consumes significant energy, and healthcare systems are increasingly pressured to reduce their carbon footprint. The added overhead of differential privacy and anomaly detection increases energy consumption per training iteration. To mitigate this, the system could adopt sparsity-inducing techniques or early stopping criteria based on prototype stability. Policy implications extend to the standardization of benchmarks for evaluating privacy-preserving healthcare AI. Currently, there is no widely accepted benchmark that simultaneously measures privacy leakage, model accuracy, and resistance to data poisoning. The research community, in collaboration with regulatory bodies, should develop a common evaluation suite that includes vertical split learning scenarios with prototype-based defenses. Such a benchmark would accelerate adoption and foster trust among stakeholders.

## **8. Future Directions and Sustainability**

The evolution of vertical split learning in healthcare will be shaped by advances in several areas. One promising direction is the development of adaptive privacy budgets that adjust based on the sensitivity of the data modality. For example, genomic data may require stricter privacy than radiology images. Prototype-based anomaly detection can be extended to a multi-class and multi-prototype setting, where each class is represented by several prototypes to capture intra-class variation. Additionally, combining prototype-based detection with generative models that reconstruct anomalies could provide explainability for clinicians. The sustainability of the system depends on its ability to scale to hundreds of hospitals without incurring prohibitive communication costs. Hierarchical architectures, where hospitals are grouped into regional clusters that each have a local coordinator, could reduce latency while preserving the benefits of vertical split learning. Energy efficiency can be improved by using quantization and pruning of the latent representations, although these techniques must be evaluated for their impact on anomaly detection performance. Another important future direction is the integration of on-device learning for wearable and edge healthcare devices. Vertical split learning can enable a smartphone app to use a small local model while sending latent features to a hospital server, but the privacy and security implications for personal devices are different from institutional settings. Prototype-based defense can still be applied, but the coordinator must handle potentially unreliable or intermittent connections. Finally, the policy landscape must evolve to provide clear guidelines on the acceptable level of privacy risk and the required robustness against adversarial attacks in medical AI. As prototype-based methods mature, they could be required as part of a certification process for AI systems used in clinical decision support.

## **9. Conclusion**

This paper has presented a comprehensive system-level analysis of a privacy-preserving vertical split learning framework for healthcare AI, augmented with prototype-based anomaly detection to defend against data poisoning. We have discussed the architectural design, privacy-enhancing mechanisms, and the integration of prototype consistency checks as a robust defense. The analysis highlighted structural trade-offs among privacy, utility, communication efficiency, and fairness, and we examined deployment challenges, governance requirements, and policy implications. The framework leverages differential

privacy and secure aggregation to protect patient data while enabling collaborative model training, and the prototype-based module provides a computationally efficient and interpretable method to detect and exclude poisoned samples. The recent work of Shui et al. on ProtoGuard-SL demonstrates the practical viability of this approach [15]. As healthcare AI becomes more pervasive, systems that balance privacy, security, and robustness will be essential. Future research should focus on adaptive privacy budgeting, scalable hierarchical architectures, energy-efficient implementations, and the development of standardized benchmarks to ensure trustworthiness across diverse clinical settings.

## References

1. McMahan, B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS).
2. Vepakomma, P., Gupta, O., Swedish, T., & Raskar, R. (2018). Split learning for health: Distributed private deep learning without sharing raw data. arXiv preprint arXiv:1812.00564.
3. Gupta, O., & Raskar, R. (2018). Distributed learning of deep neural networks via independent subnet training. arXiv preprint arXiv:1810.04521.
4. Dwork, C. (2006). Differential privacy. In International Colloquium on Automata, Languages, and Programming (ICALP).
5. Biggio, B., Nelson, B., & Laskov, P. (2012). Poisoning attacks against support vector machines. In Proceedings of the 29th International Conference on Machine Learning (ICML).
6. Gu, T., Dolan-Gavitt, B., & Garg, S. (2017). BadNets: Identifying vulnerabilities in the machine learning model supply chain. arXiv preprint arXiv:1708.06733.
7. Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. ACM Computing Surveys, 41(3), 1-58.
8. Snell, J., Swersky, K., & Zemel, R. (2017). Prototypical networks for few-shot learning. In Advances in Neural Information Processing Systems (NeurIPS).
9. Topol, E. J. (2019). High-performance medicine: The convergence of human and artificial intelligence. Nature Medicine, 25(1), 44-56.
10. Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H. R., Albarqouni, S., ... & Cardoso, M. J. (2020). The future of digital health with federated learning. NPJ Digital Medicine, 3(1), 1-7.
11. Li, T., Sahu, A. K., Talwalkar, A., & Smith, V. (2020). Federated learning: Challenges, methods, and future directions. IEEE Signal Processing Magazine, 37(3), 50-60.
12. Sweeney, L. (2002). k-anonymity: A model for protecting privacy. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, 10(5), 557-570.
13. Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., & Zhang, L. (2016). Deep learning with differential privacy. In Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS).

14. Papernot, N., Abadi, M., Ulfar, E., Goodfellow, I., Talwar, K., & Zhang, L. (2017). Semi-supervised knowledge transfer for deep learning from private training data. In Proceedings of the 5th International Conference on Learning Representations (ICLR).
15. Shui, Y., Jin, R., Dou, Z., & Gao, Z. (2026). ProtoGuard-SL: Prototype Consistency Based Backdoor Defense for Vertical Split Learning. arXiv preprint arXiv:2604.03595.
16. Bagdasaryan, E., & Shmatikov, V. (2020). Blind backdoors in deep learning models. In Proceedings of the 29th USENIX Security Symposium.
17. Sun, Z., & Li, B. (2022). A survey on data poisoning attacks and defenses in federated learning. IEEE Transactions on Neural Networks and Learning Systems, 33(11), 6073-6093.
18. Chen, Y., Qin, X., Wang, J., Yu, C., & Gao, W. (2020). FedHealth: A federated transfer learning framework for wearable healthcare. IEEE Intelligent Systems, 35(4), 83-93.
19. Zhang, C., Xie, Y., Bai, H., Yu, B., Li, W., & Gao, Y. (2021). A survey on federated learning. Knowledge-Based Systems, 216, 106775.
20. Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., ... & Zhao, S. (2021). Advances and open problems in federated learning. Foundations and Trends in Machine Learning, 14(1-2), 1-210.