# SoK: Privacy-aware LLM in Healthcare: Threat Model, Privacy Techniques, Challenges and Recommendations

Mohoshin Ara Tahera*, Karamveer Singh Sidhu†, Shuvalaxmi Dass*, Sajal Saha†

*University of Louisiana at Lafayette, Lafayette, LA, USA
†University of Northern British Columbia, Canada
Emails: mohoshin-ara.tahera1,shuvalaxmi.dass@louisiana.edu
ksidhu,sajal.saha@unbc.ca

*Abstract*—**Large Language Models (LLMs) are increasingly adopted in healthcare to support clinical decision-making, summarize electronic health records (EHRs), and enhance patient care. However, this integration introduces significant privacy and security challenges, driven by the sensitivity of clinical data and the high-stakes nature of medical workflows. These risks become even more pronounced across heterogeneous deployment environments, ranging from small on-premise hospital systems to regional health networks, each with unique resource limitations and regulatory demands. This Systematization of Knowledge (SoK) examines the evolving threat landscape across the three core LLM phases: Data preprocessing, Fine-tuning, and Inference within realistic healthcare settings. We present a detailed threat model that characterizes adversaries, capabilities, and attack surfaces at each phase, and we systematize how existing privacy-preserving techniques (PPTs) attempt to mitigate these vulnerabilities. While existing defenses show promise, our analysis identifies persistent limitations in securing sensitive clinical data across diverse operational tiers. We conclude with phase-aware recommendations and future research directions aimed at strengthening privacy guarantees for LLMs in regulated environments. This work provides a foundation for understanding the intersection of LLMs, threats, and privacy in healthcare, offering a roadmap toward more robust and clinically trustworthy AI systems.**

*Index Terms*—**LLM, healthcare, Privacy-preserving techniques, Threat model**

## 1. Introduction

LLMs are increasingly integrated into healthcare for clinical documentation, decision support, radiology and pathology report summarization, and clinician–patient communication (representative applications in Table 1). Deployments span radiology/report summarization [1], [2], triage and clinical decision support (CDS) systems [3], [4], and multilingual dialogue assistants [5], [6]. While these applications demonstrate substantial utility, they also introduce significant privacy risks due to the sensitivity of Protected Health Information (PHI) and stringent regulatory constraints.

This SoK examines privacy concerns across the three core operational phases of LLMs: data preprocessing, fine-tuning, and inference, with a specific focus on text-centric healthcare applications (e.g., EHR notes, discharge summaries, DICOM-derived reports, pathology and radiology narratives, and clinician–patient transcripts). We intentionally scope the SoK to text-generating or text-consuming LLM pipelines, excluding imaging-only models unless they interact with PHI-bearing textual artifacts.

Unlike prior surveys [7]–[12], which primarily catalog techniques or discuss privacy at a high level, our work provides a phase-aligned systematization. We unify terminology and explicitly map: attack surfaces→enabled attacks→defenses→remaining limitations grounded in healthcare data artifacts such as EHR/HL7/FHIR structures, DICOM text fields, and clinical transcripts. Prior work does not systematically link adversary capabilities to specific vulnerabilities across the three phases; Table 2 highlights these gaps. Our corpus spans peer-reviewed studies and authoritative preprints from 2020–2025 involving clinical datasets, hospital deployments, or consortium-based learning; details appear in section 2.

This SoK answers three guiding research questions in each phase:

- **RQ1.** How should adversaries in healthcare LLMs be categorized (e.g., internal vs. external), and how do their capabilities and prior knowledge shape the attack vectors that emerge at each phase of the LLM lifecycle?
- **RQ2.** How effectively do current privacy-preserving techniques mitigate phase-specific vulnerabilities?
- **RQ3.** What limitations remain in current defenses, and what phase-aware strategies can guide future privacy-enhanced LLM development?

To address **RQ1**, we construct a detailed threat model for each phase, identifying adversary capabilities, prior knowledge, and corresponding vulnerabilities. For **RQ2**, we evaluate privacy-preserving techniques and analyze their effectiveness relative to the identified phase-specific threats. For **RQ3**, we synthesize limitations and propose future directions and recommendations for phase-aware, threat-resilient, healthcare-specific privacy enhancements.

Given that healthcare deployments commonly rely on distributed architectures, our analysis of the fine-tuning stage adopts a Federated Learning (FL) framework, emphasizing vulnerabilities and defenses in the Federated Fine-Tuning Phase. By answering these questions, we make the following three contributions:

1) **Phase-Specific Threat Model:** We develop a comprehensive threat model for healthcare LLMs, catego-

rizing adversaries by location (internal vs. external) and capability and identifying attack surfaces and key vulnerabilities across data preprocessing, federated fine-tuning, and inference (e.g., gradient leakage, internal update exposure, model extraction).

2) **Evaluation of Privacy-Preserving Techniques:** We systematically analyze defenses such as differential privacy, secure aggregation, inference-time mitigations and evaluate how effectively they address the vulnerabilities surfaced in each phase. Our findings show that existing techniques provide partial protection but are often phase-agnostic or misaligned with healthcare-specific threats.

3) **Limitations and Future Directions:** We identify critical limitations in current privacy strategies and propose future research directions aimed at developing phase-aware, threat-resilient, and privacy-enhanced LLMs tailored for healthcare environments. These include recommendations such as standardized data anonymization, adaptive differential privacy to mitigate gradient leakage while preserving rare-disease fidelity in federated fine-tuning.

Following the literature collection methodology (Section 2), our SoK is structured into three sections, each focused on a distinct phase of the LLM lifecycle: Data Preprocessing (Section 3), Federated Fine-Tuning (Section 4), and Inference (Section 5). Each section includes two subsections: The Threat Model subsection addresses **RQ1**, outlining adversaries, capabilities, and attack vectors. The Privacy Preserving Defenses subsection covers the rest of the RQs: **RQ2** with *Takeaways* on how existing privacy-preserving techniques mitigate vulnerabilities, and **RQ3** with *Limitations* of current defenses and *Recommendations*. To synthesize insights across phases, we present Tables 3, 4 and 5, which systematically maps the threats and vulnerabilities of each phase to their corresponding defenses, limitations, and recommendations, offering a concise, phase-specific roadmap for privacy-aware LLMs.

## 2. Literature Collection

This SoK follows a structured but targeted methodology aimed at systematizing privacy and security risks for healthcare LLMs across the three operational phases: data preprocessing, federated fine-tuning, and inference, prioritizing works that directly engage with privacy, security, or compliance rather than providing exhaustive coverage of all clinical LLM applications. We searched IEEE Xplore, ACM DL, SpringerLink, ScienceDirect, PubMed, and major preprint repositories (arXiv, medRxiv) for the years 2019–2025 using combined terms spanning model types *("LLM," "large language model," "foundation model")*, healthcare domains *("medical," "clinical," "EHR," "radiology," "pathology," "telemedicine")*, and privacy/security mechanisms *("privacy," "differential privacy," "federated learning," "secure aggregation," "homomorphic encryption," "access control")*. Searches were supplemented with backward and forward snowballing from influential surveys on LLM privacy and healthcare AI. We included papers that (i) involve LLMs or closely related foundation models or core privacy-preserving mechanisms relevant to healthcare (e.g., DP, FL, HE/MPC, anonymization, access control); (ii) use or explicitly target healthcare data

or clinical workflows such as EHRs, imaging, pathology, telemedicine, or medical QA; and (iii) provide a substantive treatment of privacy, security, or regulatory constraints, including threat models, attacks, defenses, or evaluations. We excluded purely application-focused clinical LLM papers with only passing references to privacy, as well as generic privacy/security works lacking a clear connection to healthcare or LLMs.

## 3. Data Preprocessing Phase

The preprocessing phase involves cleaning, structuring, and transforming raw clinical data to address issues like missing values, feature normalization, class imbalance, and skewed distributions [13], [14]. In healthcare, this phase is critical due to the sensitivity and diversity of data from EHRs, medical imaging, claims, and patient-generated sources.

### 3.1. Threat Model

The data preprocessing stage ingests heterogeneous clinical data including EHR tables, radiology metadata, lab results, claims records, and patient-generated content—where identifiers and quasi-identifiers remain intact [34], [35]. Because adversaries interact with the raw substrate from which LLM datasets are formed, this phase presents uniquely powerful opportunities for privacy compromise and data manipulation [36], [37].

**3.1.1. Attacker Landscape and Incentives.** Preprocessing threats stem from both internal and external adversaries, whose incentives and privileges shape distinct privacy risks.

**Internal adversaries** are the most operationally impactful, as they handle raw data throughout daily workflows. **(1)** Clinical staff (nurses, residents, and attending physicians) may export EHR data for shift handovers or audits through informal channels, unintentionally exposing unmasked identifiers and clinical notes [35]. **(2)** Data engineers and ML researchers maintaining ETL and data-cleaning pipelines often retain identifiers for debugging or convenience, propagating PHI leakage into intermediate tables and logs [38]–[41]. **(3)** System administrators and IT personnel possess privileged access to storage servers, backups, and preprocessing nodes. Misconfigured access controls or network shares can silently expose PHI to unauthorized internal users or external exploits [39].

**External adversaries** target the same infrastructures but with financial, strategic, or competitive motives. **(1)** Ransomware groups exploit unpatched or misconfigured ETL servers and data lakes, as illustrated by the WannaCry attack on the NHS [42]. **(2)** Corporate or state-linked attackers infiltrate healthcare data lakes to harvest PHI for analytics or commercial leverage [9], [36]. **(3)** Research competitors may compromise preprocessing APIs or storage connectors to reconstruct cohorts or infer institutional practices [43]–[45].

**3.1.2. Prior Knowledge and Capability Gradient.** Adversaries in this phase often hold strong prior knowledge of healthcare data systems, making even limited access highly dangerous, such as: **(1)** They understand

TABLE 1: Representative Healthcare Applications of LLMs with example tasks

| Application Type | Example Tasks |
|---|---|
| Electronic Health Records (EHR) | Demographics, diagnoses, labs, ICU notes & pathology (MIMIC, n2c2) [13]–[15]; insurance claims [3]; rare cohorts [16]; synthetic EHRs for comorbidity/survival modeling [17]–[19] |
| Clinical Summaries & Documentation | Discharge-note assistants [4], [20]; cardiac/oncology summaries [16], [21]; multilingual notes [5], [22]; ICU timelines [23]; leakage risks (telemedicine) [24], [25] |
| Medical Imaging & Signal Data | Radiology summarization (CT, X-ray) [1], [2]; pathology summarizers [6], [21]; ChestX-ray14 [8]; cardiology data (ECG, BP, cholesterol, EF) [2], [20] |
| Clinical Decision Support (CDS) | Rule-based CDS (guidelines) [26]; oncology classification [21], [27]; cardiovascular outcome prediction [3], [4]; comorbidity prediction (graph prompting, RAG) [19]; multimodal support [27], [28] |
| Telemedicine & Patient Interaction | Triage chatbots [29], [30]; telemedicine dialogues [24]; multilingual transcripts [5], [6]; DSS for rare tumor queries [31], [32]; mobile health apps (chest pain triage) [20], [33] |

TABLE 2: Comparison of healthcare-focused surveys vs. our SoK. Legend: ★= strong/unique/comprehensive; ✓= present/addressed; ◐ = Partial; ○ = Narrow/absent; Full = Pre-Processing → Federated Fine-Tuning → Inference

| Criteria | [7] | [8] | [9] | [10] | [11] | [12] | Our SoK |
|---|---|---|---|---|---|---|---|
| Stages Covered | ◐ Partial | ◐ Partial | ◐ Partial | ◐ Partial | ◐ Partial | ◐ Partial | ★Full |
| Threat Model | ○ Narrow | ○ Narrow | ○ Narrow | ○ Narrow | ◐ Partial | ◐ Partial | ★Comprehensive, phase-mapped |
| Defenses / PPTs | ✓Partial | ✓Partial | ○ Narrow | ✓Partial | ✓Partial | ✓Partial | ★Tailored + Phase-mapped |
| Limits of Defenses | ✓Present | ○ Absent | ○ Absent | ○ Absent | ✓Present | ✓Present | ★Phase-tied |
| Recommendations | ✓Present | ✓Present | ✓Present | ○ Absent | ✓Present | ✓Present | ★Gap-tied, phase-specific |

EHR schemas and identifiers (HL7/FHIR fields, MRNs, timestamps, and ICD codes), enabling precise targeting of PHI within raw exports [13], [14], [35]. **(2)** They exploit quasi-identifiers such as age, location, ethnicity, or rare conditions for re-identification and linkage before anonymization [34], [46]. **(3)** They are familiar with ETL workflows and data handling practices, including where staging files, failed job outputs, and temporary exports are stored [38], [39]. **(4)** Such insider knowledge directly enables exploitation of weakly protected APIs and misconfigured file systems that connect preprocessing tools to hospital databases. **(5)** High-resource externals also combine auxiliary datasets such as insurance claims, leaked registries, and prior hospital breaches to strengthen re-identification or model poisoning efforts [47].

The capability spectrum spans: **(1)** low-resource insiders (e.g., nurses, clerks) can unintentionally leak identifiers through manual exports or shift handover lists. [34], [35]; **(2)** moderate-resource insiders (e.g., ETL engineers, ML staff) handle staging data and debugging logs containing unmasked PHI. [40], [41]; and **(3)** high-resource externals (e.g., ransomware groups, APTs) exploit insecure APIs, exfiltrate backups, or inject poisoned HL7/FHIR records to compromise downstream model integrity. [48], [49].

**3.1.3. Attack Surface Vulnerabilities and Enabled Attacks.** These layered attacker capabilities map directly to the core vulnerabilities of the preprocessing stage:
- **Quasi-identifiers:** Demographic and clinical combinations exploited by insiders and externals for re-identification [34].

- **Data leakage via logs:** Debug outputs and ETL failure logs containing PHI accessible to authorized insiders or compromised systems [41].
- **Weak access controls and APIs:** Misconfigured permissions, exposed endpoints, or insecure ETL connectors exploited by administrators or external attackers to access raw datasets [38], [39].

These vulnerabilities directly enable several concrete attacks: **(1) Re-identification/linkage attacks**, where adversaries combine quasi-identifiers with auxiliary records to recover patient identities [46]. **(2) Log-based PHI leakage**, where unmasked identifiers in ETL traces or debug outputs expose raw clinical details [40]. **(3) Unauthorized dataset access**, where weak permissions or misconfigured APIs allow direct retrieval of raw EHR exports or staging files [39]. **(4) Data poisoning**, where attackers inject manipulated or fabricated clinical records into HL7/FHIR streams to distort downstream fine-tuning behavior [49].

Together, these attacks operationalize the incentives, capabilities, and prior knowledge described in the pre-processing threat model, demonstrating how seemingly routine workflow weaknesses become attack vectors.

## 3.2. Privacy-Preserving Defenses

Given the risks of handling PHI, implementing robust privacy-preserving techniques in this phase is crucial [50]. Common methods like *data anonymization*, *synthetic data generation*, and *noise addition* help safeguard patient privacy while enabling LLM training.

TABLE 3: Data Preprocessing Phase — Summary of Threats, Defenses, Limitations, and Recommendations.

| Threat Model | Defenses (PPTs) | Limitations | Recommendations |
|---|---|---|---|
| • **Adversaries:** Internal (clinicians, data engineers, IT admins) and external (ransomware, state-linked, competitors). <br> • **Capabilities:** Exploit quasi-identifiers, extract PHI from ETL logs/debug traces, abuse weak APIs, or inject poisoned HL7/FHIR data. <br> • **Key Vulnerabilities:** Quasi-identifiers, PHI-bearing logs, weak access control, raw exports/backups. | • **Anonymization:** Masking, tokenization, k-anonymity, l-diversity, pseudonym vaults. <br> • **Synthetic Data:** GAN/VAE/LLM-based generation replacing real EHR inputs. <br> • **Differential Privacy:** DP-SGD, local DP, stochastic embedding noise for pre-training or tabular features. | • **Quasi-identifiers:** Generators may reproduce unique attribute patterns, enabling linkage. <br> • **Logs / APIs:** PHI may persist in debug traces; DP misconfiguration leaks identifiers. <br> • **Poisoning:** Defenses fail against manipulated records injected pre-noise or masking. <br> • **Governance:** Weak auditing or inconsistent DP parameters reintroduce vulnerabilities. | • **Quasi-identifiers:** Structured anonymization within ETL (tokenization, suppression). <br> • **Log leakage:** LLM-assisted scrubbing of debug/staging data. <br> • **Access control:** Least-privilege, vault-based tokens and synthetic sandboxes. <br> • **Poisoning:** Pre-ingestion anomaly detection and HL7/FHIR integrity checks. |

**3.2.1. Data Anonymization.** It is the first operational defense in preprocessing, directly countering insider misuse and quasi-identifiers exposure before clinical data enter modeling pipelines. During cleaning, raw EHR tables, radiology metadata, and lab files still carry names, MRNs, timestamps, and geocodes, which are the prime targets for careless or malicious insiders [35], [38], [40], [41]. As outlined in the threat model, low-resource insiders (nurses, clerks) leak unmasked exports in handovers, and moderate-resource insiders (ETL/ML staff) surface identifiers in debug logs and staging data [39]. Early masking or tokenization interrupts these leakage paths so exports, logs, and cached views hold only pseudonymous values even when systems are misconfigured.

Classical frameworks such as k-anonymity, l-diversity, t-closeness generalize or perturb demographic/clinical fields so attribute combinations cannot uniquely identify a patient [34], [51], [52]. This directly mitigates the linkage risks noted in the threat model, where adversaries combine auxiliary attributes (age, ZIP, rare conditions) to re-identify patients [34], [46]. Hospitals commonly truncate postal codes or merge rare diagnoses (e.g., oncology registries), blunting the advantage of external actors who leverage demographic priors or leaked registries [47]. Complementary pseudonymization and encryption replace identifiers with reversible tokens kept in audited vaults [53], [54], constraining the practical privileges of ETL personnel and admins: if staging/backup files are accessed, only encrypted tokens, not PHI, are exposed, cutting off the same audit log and backup channels exploited by insider error or lateral movement [40], [41].

**Using anonymized datasets.** Most healthcare-LLM studies lean on pre-anonymized corpora (MIMIC-III/IV, MIMIC-CXR) de-identified for HIPAA/GDPR [13], [14], [53], [55], enabling clinical IE [3], privacy-preserving local deployment [56], and federated pretraining [15]. Yet dataset-level compliance does not eliminate pipeline-level exposure: identifiers can resurface in temporary exports, logs, or monitoring dashboards during hospital preprocessing [40], and quasi-identifiers can still support attribute inference or record linkage [43], [57]. Some studies instead use closed-source anonymized sets [58] or "PII-free by design" resources (ChatDoctor, MMQS, IMCS-21) [6], [59], which avoid direct identifiers but rarely provide end-to-end audit evidence that PHI never entered the pipeline.

**Using LLMs for anonymization.** An emerging mitigation embeds anonymization inside preprocessing itself:

LLMs fine-tuned for medical de-identification redact PHI while preserving clinical semantics [60]. These tools address the free-text/logging vulnerabilities highlighted in the threat model, e.g., identifiers buried in notes or timestamps that leak via ETL logs and debug outputs [38], [40] and, when paired with hybrid NER–DL pipelines, sanitize text streams before storage or model ingestion, closing insider and scraper pathways that rely on residual identifiers for linkage [47]. Implementation quality, however, is uneven: some datasets report manual redaction (e.g., HajjHealthQA) [27], others assume PHI absence in synthetic dialogues [61], [62], and several open resources still show explicit identifiers without automated scrubbing [63]. This inconsistency reinforces the need to treat anonymization as an engineered control in the preprocessing pipeline, not merely a property of upstream datasets.

**Takeaway.** Anonymization directly mitigates the preprocessing threats by disrupting how insiders and externals exploit raw clinical data before transformation. Masking or encrypting identifiers before ETL prevents *re-identification* and *linkage attacks* via quasi-identifiers such as age or rare conditions, while also blocking *log-based PHI leakage* from debug outputs [34], [40]. For *high-resource external adversaries*, it removes identifiers that enable cross-record matching through their prior knowledge and auxiliary datasets, thereby reducing the utility of stolen clinical data. Consistent application of this technique substantially reduces exposure to core preprocessing vulnerabilities: *quasi-identifiers misuse, log leakage, and weak access controls*, forming a crucial first layer of defense for privacy-preserving LLM pipelines.

**Limitation.** While anonymization mitigates key preprocessing threats, it leaves several vulnerabilities partially exposed. *(1) Quasi-identifiers:* it cannot fully prevent re-identification or linkage attacks by high-resource externals leveraging prior knowledge and auxiliary datasets; *(2) Data leakage via logs:* when applied only at the dataset level, it fails to mitigate PHI exposure through debug outputs, temporary exports, or audit traces, allowing low- and moderate-resource insiders to access identifiers left in system logs or misconfigured dashboards. *(3) Weak access controls:* anonymization does not prevent data poisoning or unauthorized input manipulation, as moderate-resource insiders can still inject or alter records before masking through insecure ETL access points.

**3.2.2. Synthetic Data Generation.** This operates as a proactive privacy defense in the preprocessing phase, directly addressing the insider misuse and external re-identification threats outlined in the threat model. Whereas anonymization masks identifiers after data collection, synthetic generation removes reliance on raw PHI altogether, eliminating exposure points before they arise. By simulating clinical distributions from learned or rule-based patterns, it dismantles the economic and operational incentives for insiders who might leak raw EHR exports and renders external ransomware operators unable to monetize stolen backups or staging snapshots [36], [48]. Because synthetic records contain no true identifiers or one-to-one mappings to real patients, they nullify quasi-identifier risks and schema-specific leak paths that low- and moderate-resource insiders can trigger during data preparation [35], [40]. For high-resource external adversaries, the value of exfiltrated data collapses; synthetic ETL snapshots or exported samples cannot be used for re-identification, blackmail, or population linkage [46], [47].

**Rule-based methods:** Rule-driven and simulation approaches generate synthetic data from explicit clinical logic or population statistics, providing auditable and deterministic privacy guarantees. Monte Carlo and discrete-event models replicate hospital workflows—admissions, lab requests, and comorbidities, without using real patient identifiers [64], [65]. For instance, the Dismed dataset randomized annotated entities using biomedical ontologies to maintain diagnostic structure while erasing sensitive identifiers [60], [66]. Such systems are particularly effective against insider leakage because their generation logic never touches PHI; even if preprocessing pipelines or exports are compromised, the data itself is synthetic and devoid of re-identifiable features.

**Using LLMs:** Learning-based methods extend this protection to complex multimodal data such as radiology reports, clinical narratives, and tabular EHRs. Generative networks (GANs, VAEs, and diffusion models) and LLMs capture nonlinear dependencies across diagnoses, medications, and outcomes [67], [68]. LLM-based generators—such as GPT-4 and LLaMA 3.1-70B, have been used to produce synthetic case studies and comorbidity graphs for models like ComLLM and Asclepius, supporting clinical QA and CDS tasks without direct patient data access [18], [19], [69]. This approach directly mitigates the text-based leakage vulnerabilities identified in the threat analysis, where identifiers persist in unstructured notes or debugging logs [38], [40]. By inserting a synthetic layer between hospital data and LLM pipelines, healthcare institutions convert high-risk preprocessing workflows into auditable, privacy-preserving sandboxes.

**Takeaway.** Synthetic data generation directly mitigates the vulnerabilities identified in the preprocessing threat model by removing dependence on *raw PHI* before training. By replacing true records with statistically valid but fictitious samples, it disrupts the *economic and operational incentives* of low- and moderate-resource insiders who rely on *understanding of EHR schemas and identifiers* and *access to staging data* to exploit *quasi-identifiers* or export unmasked EHRs. At the same time, synthetic data eliminates the informational value that high-resource external adversaries derive from *auxiliary datasets*, rendering *re-identification*, *linkage attacks*, and exploitation

of *EHR exports* or *staging leaks* ineffective. By removing the real identifiers that make these vulnerabilities actionable, synthetic data closes the same *attack surfaces* described in the threat model and reshapes the preprocessing pipeline into a controlled environment where both insider misuse and external infiltration have substantially reduced impact.

**Limitation.** While synthetic data generation reduces dependence on real PHI, it does not eliminate preprocessing vulnerabilities. *(1) Quasi-identifiers:* generators may inadvertently reproduce statistical or quasi-identifier patterns, leaving the system exposed to the same re-identification and linkage risks exploited by high-resource adversaries with strong prior knowledge and auxiliary records. *(2) Weak access controls:* synthetic pipelines remain vulnerable to data poisoning introduced before generation, allowing moderate-resource insiders or malicious collaborators to manipulate raw PHI and exploit their access privileges. *(3) Data leakage via logs:* if generation quality is inconsistent or insufficiently audited, synthetic datasets may recreate structural cues that allow adversaries to exploit previously identified attack surfaces such as EHR exports, staging leaks, or API exposures. Thus, while synthetic generation reduces reliance on real PHI, its effectiveness depends on strict governance and validation to prevent threat-model vulnerabilities from re-emerging through statistical artifacts.

**3.2.3. Differential Privacy (DP).** provides a formal defense against the inference and memorization risks that persist in preprocessing, particularly when partial identifiers or structured embeddings are exposed. Unlike anonymization or synthetic generation, DP offers quantifiable privacy guarantees by injecting calibrated noise into data values, gradients, or outputs [70]. This directly addresses the moderate-resource insiders and external attackers identified in the threat model, those capable of inspecting gradients, intermediate logs, or fine-tuning checkpoints [38], [40]. By applying Gaussian or Laplacian noise during aggregation or embedding generation, DP ensures that individual patient records remain statistically indistinguishable, even to privileged ETL engineers or model developers [71], [72]. Healthcare frameworks such as DisLLM, MedMCQ, and PubMedQA demonstrate how DP-SGD and local DP protect sensitive EHR and clinical text embeddings while retaining acceptable task accuracy [72]–[74]. Local DP, used in resource-limited hospital preprocessing, adds noise directly at the feature extraction stage, obscuring PHI before it reaches shared systems, while centralized DP reduces cross-phase linkage between preprocessing and downstream fine-tuning [8]. Beyond strict DP, stochastic embedding methods such as NEFtune and SHADE-AD [23], [75] introduce non-determinism to reduce overfitting and memorization. Although lacking formal guarantees, these strategies complement DP by further diminishing the information advantage of internal or external adversaries during preprocessing.

**Takeaway.** DP directly mitigates the inference and reconstruction risks described in the preprocessing threat model by protecting against the misuse of *gradients*, *intermediate logs*, and *structured embeddings*. By adding calibrated noise during feature extraction or aggregation, DP reduces the ability of moderate-resource insiders, who

have access to *debugging logs*, *staging data*, or *model checkpoints*, to recover patient attributes or exploit *quasi-identifiers*. The same noise also weakens external adversaries who rely on auxiliary datasets to perform *re-identification* or *record reconstruction*. In this way, DP converts preprocessing from a trust-dependent stage into one that enforces quantifiable privacy guarantees, directly addressing the vulnerabilities associated with linkage attacks and log leakage identified in the threat model.

**Limitation.** Despite its benefits, DP does not address all vulnerabilities in the preprocessing threat model. *(1) Weak access controls/APIs:* DP cannot stop data poisoning, since adversaries can inject or manipulate raw clinical records before noise is applied, exploiting insecure ETL workflows. *(2) Quasi-identifiers:* the privacy–utility trade-off can distort sensitive fields, degrading fidelity in downstream clinical tasks while leaving partial re-identification risks. *(3) Data leakage via logs:* DP's effectiveness depends on correct configuration; mis-set noise budgets or implementation flaws can still expose PHI through debug traces or unsecured system logs. Thus, while DP mitigates re-identification and log exposure, it cannot fully prevent poisoning or misconfiguration exploits in preprocessing.

---

**Recommendation 1:** The preprocessing phase exposes the most sensitive attack surface of the healthcare LLM pipeline, where adversaries exploit raw identifiers, staging artifacts, and weak system configurations. To mitigate these vulnerabilities, defenses must directly align with the threat model. **(1) Quasi-identifiers Exposure.** Implement task-specific anonymization combined with differential privacy audits to remove or mask high-risk identifiers before ETL execution. This reduces re-identification and linkage attacks that exploit demographic and clinical quasi-identifiers using prior knowledge and auxiliary datasets **(2) Data leakage via logs.** Adopt privacy-aware logging frameworks with automatic redaction, hashed identifiers, and secure audit retention policies. This ensures ETL traces and debugging outputs cannot expose PHI through intermediate tables or system logs accessible to insiders or compromised systems. **(3) Weak access controls and API exposures.** Deploy role-based access control (RBAC) and zero-trust authentication across preprocessing nodes, ensuring that only verified users and processes can query or export data. Integrate secure API gateways and encryption-in-use (TEE) mechanisms to protect intermediate exports and prevent unauthorized dataset retrieval **(4) Data poisoning.** Incorporate data integrity validation and schema-constrained ingestion to detect anomalous or manipulated records injected before anonymization. Combine this with input provenance tracking and hash-based cross-validation to ensure the authenticity of clinical feeds (HL7/FHIR) before preprocessing pipelines ingest

---

Table 3 summarizes the preprocessing-phase taxonomy, consolidating internal/external actors, attack surfaces, defenses, limitations, and actionable recommendations derived from this section.

# 4. Federated Fine-Tuning Phase

Building on the vulnerabilities introduced during preprocessing, the fine-tuning phase exposes a new class of distributed risks: even when raw PHI is masked, residual artifacts, cohort structures, and poisoned records propagate into gradients and client updates across federated hospitals. We focus specifically on federated fine-tuning because, in healthcare, centralized fine-tuning is rarely feasible in cross-institutional raw-data pooling violates HIPAA/GDPR, institutional policies, and multi-site data-sharing agreements [8], [9], [34], [35]. As a result, federated fine-tuning is not an architectural preference but the only legally and operationally viable mechanism for adapting LLMs to clinical data in real deployments [4], [10], [15], [76]. This makes FL the principal attack surface and hence the appropriate scope for a phase-aware SoK.

Fine-tuning LLMs in healthcare enables models to adapt to domain-specific terminology and tasks, improving accuracy and clinical relevance. However, fine-tuning on sensitive, institution-specific datasets introduces distinct privacy risks. While techniques such as Differential Privacy (DP) and adapter tuning offer partial protection [77], the majority of practical deployments and research indicate that Federated Learning (FL) is the most comprehensive privacy-preserving framework for regulated healthcare environments [4], [21], [78]. In this phase, we outline the threat landscape specific to FL, focusing on the privacy challenges arising from distributed training infrastructures connecting hospitals.

## 4.1. Threat Model

In the fine-tuning phase, threats arise from adversaries, internal or external, who interact with the FL infrastructure connecting hospitals. Unlike the preprocessing stage, where unmasked raw data is exposed, here the primary privacy risks emerge from access to *local gradients*, *client updates*, and *communication buffers*. These artifacts implicitly encode sensitive clinical attributes and thus create a distinct attack surface.

**Internal adversaries** pose the most operationally impactful risks because they operate inside trusted hospital boundaries and interact directly with local training workflows: **(1)** Clinicians or annotators may inadvertently expose local patient records while validating model predictions or providing feedback on draft outputs [2]. **(2)** ML engineers and data scientists routinely inspect gradient snapshots or checkpoints during debugging, unintentionally accessing encoded EHR information [20]. **(3)** System administrators and IT staff managing synchronization or storage systems have privileged visibility into local model caches, communication buffers, or server logs, making silent extraction or export of update data feasible [21]. These internal actions rarely trigger alarms because they occur within expected operational workflows.

**External adversaries** target the distributed nature of FL and the heterogeneity of hospital networks: **(1)** Network intruders intercept parameter updates or partial gradients in transit, enabling reconstruction of patient details via gradient inversion [43]. **(2)** Corporate or state-linked actors inject malicious updates to bias shared clinical models, impairing diagnostic consistency across sites [5].

**(3)** Ransomware and extortion groups target aggregation servers to exfiltrate multi-institutional model checkpoints encoding sensitive patterns across cohorts [4]. These adversaries are driven by financial, strategic, or competitive motives and exploit FL's distributed communication as a large-scale leakage vector.

### 4.1.1. Prior Knowledge and Capability Gradient. Adversaries in this phase vary in both domain knowledge and computational capacity, which determine how effectively they can exploit gradients, checkpoints, or update buffers: **(1)** Low-resource insiders (e.g., clinicians, annotators) understand how patient-level features are represented in model behavior and may unintentionally expose encoded EHR details during validation or feedback. **(2)** Moderate-resource insiders (data scientists, ML engineers) possess deeper knowledge of model architectures and storage layouts, knowing where gradient files, synchronization buffers, or checkpoints are stored, allowing silent access or export during debugging or maintenance [21]. **(3)** High-resource external adversaries, such as coordinated ransomware groups or corporate actors, leverage prior knowledge of biomedical embeddings and auxiliary EHR datasets to align intercepted gradients with public checkpoints, reconstructing rare conditions or site-specific vocabulary [43], [70].

This knowledge gradient corresponds directly to their operational capabilities: **(1)** Low-resource insiders leak updates through logs and validation workflows; **(2)** Moderate-resource insiders or externals perform gradient inversion and membership inference using intercepted traffic or open biomedical corpora [5], [79]; and **(3)** High-resource actors conduct model alignment or poisoning across hospitals, exploiting the distributed update-sharing and aggregation layers [4], [20]. Together, these escalating capabilities explain how prior knowledge, from local system familiarity to cross-institutional data access, translates directly into exploitation of fine-tuning vulnerabilities such as gradient leakage, client update interception, and poisoning within federated healthcare networks.

### 4.1.2. Attack Surface Vulnerabilities and Enabled Attacks. Distinct attack surface vulnerabilities emerge during fine-tuning of healthcare LLMs, each tied to the nature of clinical datasets and federated infrastructures.

- **Model gradient leakage:** During fine-tuning, gradients or model weights stored locally may inadvertently memorize sensitive clinical data such as structured EHR variables (e.g., blood pressure, cholesterol levels, ICD-10 codes) and unstructured discharge narratives. Even partial leakage can expose individual patient trajectories, particularly for rare diseases recorded in small data sets. For example, in cardiology, the leakage of gradients can reveal critical patient outcomes like ejection fraction or surgical interventions—easily re-identifiable in small cohorts [16], [20]

- **Client update leakage:** In federated fine-tuning, hospitals share model updates instead of raw client data. However, these updates still carry sensitive EHR patterns, such as lab panels, clinical notes, or medication histories, which can still be leaked. Adversaries may use gradient inversion to reconstruct detailed narratives (e.g., ICU notes, pathology reports). Moreover, poisoning attacks may skew diagnoses, such as for diabetes or heart failure [4], [15]. This makes federated updates particularly sensitive, as partial feature leakage can compromise longitudinal EHR timelines.

- **Communication channel interception:** During model synchronization, adversaries who intercept updates can recover latent embeddings associated with sensitive clinical data such as EHR-derived phenotypes, lab trajectories for oncology patients, or cardiology monitoring logs. In multilingual hospital systems, intercepted updates can also expose doctor–patient transcripts integrated with EHRs, potentially revealing sensitive lifestyle or family history details tied to EHRs [5], [76].

- **Misconfigured audit and logging systems:** Many logging mechanisms often capture sensitive traces present in model training data, including PHI tokens from EHR fields such as medication lists, allergies, or comorbidities. If these logs are misconfigured or improperly secured, adversaries can gain access to sensitive metadata such as genetic risk markers or rare disease cohorts, which could then be cross-linked with external datasets [3], [6].

These vulnerabilities directly enable several concrete attack types in federated fine-tuning: **(1) Gradient inversion attacks**, where adversaries reconstruct sensitive clinical details, such as ICU notes, pathology descriptions, or lab trajectories from leaked or intercepted gradients. **(2) Membership inference attacks**, which allow adversaries to determine whether a specific patient's records contributed to the fine-tuning process, particularly for rare diseases or small-site cohorts. **(3) Update leakage attacks**, where model updates shared across hospitals reveal structured or unstructured EHR patterns embedded in weight deltas or optimizer states. **(4) Model poisoning and backdoor attacks**, where malicious insiders or externals inject crafted updates that bias diagnosis-related outputs or embed hidden triggers into clinical prediction tasks. **(5) Communication-layer interception attacks**, where man-in-the-middle adversaries extract latent representations from synchronization streams exchanged between hospitals and the central aggregator.

## 4.2. Privacy-Preserving Defenses

This section reviews core privacy defenses integrated into federated learning (FL) fine-tuning workflows covering client-side, secure update sharing, and communication safeguards followed by limitations and recommendations.

### 4.2.1. Client Side. These cover the privacy-preserving techniques employed at the client side during training in the FL setup.

**Differential Privacy (DP)** adds calibrated noise to model updates before leaving the hospital, preventing patient-level re-identification in sensitive datasets such as MIMIC-IV ICU notes or oncology records [8]. DP-LoRA applies noise only to adapter layers, preserving accuracy [80], while selective DP targets the most sensitive parameters [78], [81]. Despite potential signal loss in rare disease cohorts [16], DP remains a practical safeguard

for HIPAA/GDPR compliance and multi-hospital training [82].

**Secure Multi-Party Computation (SMPC)** enables hospitals to train jointly without sharing raw data. Each site encrypts and splits its updates, which are only usable when aggregated [82]. This supports legally restricted studies like cancer survival or cardiovascular prediction [3]. While secure against reconstruction, SMPC's high compute cost limits real-time applications [8], though it remains vital for regulated hospital consortia [4].

**Split Learning (SL)** partitions the model between client and server—local layers process raw data, and only intermediate features are shared [8], [82]. This protects imaging and textual data (e.g., ChestX-ray14, ECG datasets) while supporting multimodal tasks. Though intermediate features may leak partial identity traces, SL's lightweight setup enables participation from smaller hospitals with limited infrastructure.

**Randomized Low-Rank Adaptation (LoRA)** fine-tunes only low-rank parameters with added randomization, hindering inversion attacks on sensitive text (e.g., oncology notes, psychiatric transcripts) [80], [81]. It offers better accuracy than DP and can be combined with it for stronger protection [28]. Its low compute footprint makes it ideal for hospitals with limited GPUs, balancing privacy, efficiency, and performance.

**Quantization** reduces weight precision (e.g., 32-bit $\rightarrow$ 8/4-bit), limiting exploitable detail while cutting memory and bandwidth needs [3], [82]. Effective on structured medical data, quantized models maintain accuracy and resist gradient inversion [83]. When paired with DP or LoRA, it forms a layered defense which is resource-efficient and scalable for diverse healthcare networks [4], [78], [84].

**Takeaway.** Client-side mechanisms directly address privacy risks arising from local training gradients and intermediate checkpoints. *DP* mitigates memorization and gradient inversion by adding calibrated noise, reducing the ability of internal or external adversaries to reconstruct structured or narrative EHR data from local model states. *SMPC* encrypts local gradients before aggregation, preventing adversaries from extracting embeddings or conversational features from intercepted updates. SL minimizes raw data transmission by keeping sensitive features (e.g., imaging, ECG traces, pathology text) local and sharing only intermediate activations. *Randomized LoRA* introduces stochasticity in parameter updates, weakening the consistency needed for gradient correlation or poisoning. *Quantization* further obscures precise gradient values, reducing leakage in logs and audit traces while lowering communication overheads. Together, these defenses collectively address gradient leakage, poisoning, and metadata exposure across distributed clients.

**4.2.2. Client Update Sharing and Secure Aggregation.** This stage targets vulnerabilities in update transmission and aggregation, where even privacy-preserving local training can expose sensitive patterns. Because the aggregator is a central risk point, mechanisms such as secure aggregation and blockchain-based FL are critical in healthcare to maintain both confidentiality and institutional accountability.

**Secure Aggregation** ensures that hospital updates remain confidential during cross-silo fine-tuning. Each client masks its gradients or adapter updates with random values; only the summed result reveals the true aggregate [4], [43], [81]. This prevents reconstruction of sensitive features such as medication histories or lab trajectories. To reduce overhead, recent work masks only LoRA adapter weights instead of full gradients [78], [80], while anomaly detection helps flag poisoned updates targeting rare-disease cohorts [3]. Though computationally demanding for small hospitals [82], secure aggregation remains essential in federated healthcare networks, balancing privacy, compliance, and trust.

**Weight Delta Sharing.** transmits only parameter differences between local and global models, minimizing exposure of raw patient data [83]. This approach improves efficiency for EHR-based LLMs such as MIMIC-IV discharge models [23]. LoRA-only delta sharing compresses updates into low-rank matrices for multilingual hospital networks [5], while selective layer freezing focuses on clinically relevant upper layers [16], [21]. However, repeated deltas can leak sensitive trends; combining delta sharing with DP or gradient clipping mitigates this risk [81].

**Blockchain for Update Integrity and Unlearning.** provides traceability and auditability by recording encrypted update hashes, timestamps, and institutional IDs [79]. This deters tampering and enables federated unlearning, removing a client's contributions without full retraining [28], [79]. Provenance records further enhance institutional trust in federated healthcare deployments [82]. While added infrastructure and latency pose challenges, blockchain remains valuable for high-stakes domains such as oncology or ICU consortia [21], [78].

**Takeaway.** At the aggregation layer, update-sharing mechanisms mitigate vulnerabilities associated with cross-site synchronization and multi-round leakage. *Secure ag-*

TABLE 4: Federated Fine-Tuning Phase — Summary of Threats, Defenses, Limitations, and Recommendations.

| Threat Model | Privacy-preserving Defenses | Limitations | Recommendations |
|---|---|---|---|
| <ul><li>**Internal adversaries:** clinicians/annotators (validation leaks); ML engineers (gradients, checkpoints); system admins (caches, sync buffers).</li><li>**External adversaries:** network intruders, ransomware groups, state-linked or corporate actors.</li><li>**Capabilities:** gradient inversion, membership inference, poisoning/backdoors, alignment with auxiliary corpora, client-update reconstruction.</li><li>**Key Vulnerabilities:** local gradients, client updates, optimizer states, communication buffers, synchronization traffic, misconfigured logs/audits.</li></ul> | <ul><li>**Client-side:** DP-SGD, DP-LoRA, selective/local DP; SMPC for encrypted/split updates; Split Learning (local early layers); randomized LoRA; quantization for low precision (8/4-bit).</li><li>**Client update sharing:** Secure Aggregation (masking, LoRA-only); Weight-Delta sharing (adapter updates only); Blockchain for integrity, tamper-evidence, and unlearning.</li><li>**Communication channel:** Adapter-based compression (LoRA), few-shot learning (small gradients), and RTIR (lightweight reasoning fine-tuning).</li></ul> | <ul><li>**Model gradient leakage:** DP degrades rare-disease fidelity; cached tensors persist; Split Learning leaks intermediate activations.</li><li>**Client update leakage:** SMPC protects only at aggregation; poisoning and backdoors persist; deltas reconstructable across rounds.</li><li>**Communication interception:** LoRA lowers size but not local leakage; few-shot gradients reconstructable; RTIR depends on secure retrieval.</li><li>**Misconfigured logs/audits:** logs capture gradients/activations; blockchain metadata can reveal institutional identifiers.</li></ul> | <ul><li>**Gradient leakage:** use adaptive DP with randomized LoRA to reduce inversion and membership inference.</li><li>**Update leakage:** combine SMPC with clipping and anomaly detection; selectively freeze sensitive layers.</li><li>**Communication interception:** adopt authenticated encryption, quantized aggregation, and TEEs; pair LoRA with DP.</li><li>**Audit exposure:** use privacy-aware provenance; minimize logged gradients; obfuscate ledger metadata.</li><li>**Combined mitigation:** integrate quantization + adaptive DP + randomized LoRA for fidelity, low bandwidth, and poisoning resistance.</li></ul> |

*gregation* limits client update leakage at the aggregator and weakens gradient inversion and membership inference on per-site contributions by preventing inspection of individual updates. *Weight-delta sharing* narrows exposure but must be paired with DP or clipping to resist multi-round leakage and subsequent update reconstruction. *Blockchain-based provenance* strengthens protection against misconfigured audit and logging systems and helps surface poisoning/backdoors by providing tamper-evident integrity and traceability for updates. Quantized aggregation alleviates communication-channel bottlenecks by compressing updates, though aggressive compression can affect rare-disease fidelity and bias.

**Limitation.** Update-sharing defenses reduce some risks but still leave critical gaps when mapped to the attack surfaces. *(1) Client update leakage:* Secure aggregation conceals individual hospital updates at the aggregator, but it does not prevent poisoning/backdoors adversaries can inject crafted, valid-looking shares, and endpoint artifacts (caches, optimizer states) can still leak per-site information. *(2) Multi-round leakage (client update leakage):* Weight-delta sharing lowers bandwidth per round, yet deltas accumulated across rounds enable reconstruction and membership inference over pathology or ICU trajectories, sustaining the multi-round leakage surface without additional DP or clipping. *(3) Misconfigured audit and logging systems:* Blockchain provides tamper evidence, not confidentiality; if ledger or audit metadata (institution IDs, timestamps, cohort counts) are exposed or logs are unsanitized, cross-linking and profiling of rare-disease cohorts remain possible on the audit surface. *(4) Communication channel interception:* Quantized aggregation compresses updates but does not secure transport or endpoints; intercepted compressed streams can still be analyzed, and aggressive rounding degrades rare-disease fidelity, opening room for biased predictions in oncology or cardiology tasks.

**4.2.3. Communication Channel.** As LLMs integrate into federated learning (FL), communication overhead becomes a major constraint, distinct from classical FL threats. In healthcare, transmitting full model updates is impractical; hence, reducing transmission size and frequency is crucial for scalability, privacy, and performance.

**Adapter-based Compression** such as LoRA minimizes communication by updating only low-rank adapter parameters rather than full model weights [80]. This allows hospitals to share compact updates while keeping billions of base parameters frozen, cutting bandwidth use [28], [83], [85]. Studies such as Med42 confirm LoRA-based fine-tuning maintains accuracy for oncology classification and discharge summary generation while significantly reducing parameter size [16]. LoRA-only updates also enable multilingual medical transcript training across silos [5] and limit PHI exposure since fewer parameters leave the institution [10], [81]. Additionally, adapter-level sharing supports selective unlearning, removing hospital-specific contributions without retraining [21]. However, LoRA may underperform on deep multimodal tasks like Clipsyntel summarization that require richer representations [27].

**Few-shot Learning** reduces communication by fine-tuning on limited local samples (10–50), yielding sparse, lightweight gradients [76], [84]. Applied to ICU discharge notes and oncology reports, this method transmits only task-relevant updates while maintaining accuracy [23]. Few-shot FL further supports multilingual transcripts for hospitals with small datasets [5], [22]. Smaller updates inherently reduce inversion risk [81], and combining few-shot training with DP or adapter tuning enhances stability [86]. Still, performance drops in multimodal settings like Clipsyntel-based question summarization, which demands deeper contextual learning [27].

**Real-Time Information Retrieval (RTIR)** decouples knowledge from model weights, retrieving external data during inference instead of embedding it in updates [5], [78]. In healthcare, this allows querying clinical knowledge bases for oncology guidelines, multilingual symptom data, or ICU protocols without transmitting sensitive gradients. Hospitals then fine-tune only lightweight reasoning

layers, dramatically shrinking update size and preventing patient data memorization [6]. RTIR integrated with adapter tuning supports efficient and privacy-preserving federated pipelines while sustaining diagnostic accuracy [87]. The trade-off lies in reliance on secure retrieval infrastructure, which may challenge smaller hospitals, and its limited autonomy in decision-making. Nonetheless, RTIR offers a dynamic, privacy-resilient strategy for fast-evolving healthcare domains.

**Takeaway.** Channel-level defenses map to the vulnerabilities identified in the threat model. Adapter-based Compression (LoRA) reduces *bandwidth usage* and narrows client update leakage by transmitting only adapter updates, weakening opportunities for gradient inversion and membership inference over communication buffers. Few-shot Learning yields sparse, smaller updates from limited local examples, further reducing the signal available for update reconstruction on sensitive EHR trajectories and pathology timelines and lowering exposure during *communication channel interception*. RTIR decouples knowledge from model weights so entire clinical histories are not embedded in checkpoints, shrinking the footprint of client updates and the attack surface for interception and downstream leakage. Collectively, these methods target *interception, leakage, and channel vulnerabilities*, reducing adversarial opportunities in federated healthcare synchronization streams.

**Limitation.** Communication-efficient methods reduce bandwidth but still leave key threat surfaces exposed in healthcare FL. *(1) Model gradient leakage:* Adapter-based compression (LoRA) restricts parameter sharing but does not sanitize local artifacts; reduced capacity can weaken rare-disease fidelity while leaving residual signal for gradient inversion and membership inference on the gradient/memorization leakage surface. *(2) Client update leakage:* Few-shot learning minimizes exchanged updates, yet sparse gradients can omit subtle clinical signals (e.g., drug–drug interactions) and remain vulnerable to multi-round reconstruction and poisoning/backdoors, sustaining exposure on the client update surface and risking biased predictions for under-represented tasks. *(3) Communication channel interception:* RTIR avoids embedding PHI in weights, but insecure retrieval paths and cached queries keep the communication buffers susceptible to interception and cross-linking (e.g., transcripts, sensitive queries for HIV or psychiatric histories).

> **Recommendation 2:** To close the defensive gaps identified in the threat model, future federated healthcare systems should employ layered, phase-specific defenses explicitly mapped to the key vulnerabilities. **(1) Model gradient leakage.** Adaptive DP should be combined with Randomized LoRA to mitigate gradient inversion and memorization risks while preserving fidelity for rare-disease or small-cohort training. This pairing directly strengthens protection against internal engineers or external interceptors who exploit gradient sensitivity.**(2) Client update leakage.** Augment Secure Aggregation and Weight-Delta Sharing with gradient clipping and anomaly detection to detect poisoning and reconstruction attempts before aggregation. These measures counter moderate-resource

adversaries who manipulate or infer sensitive EHR patterns from multi-round updates. **(3) Communication channel interception.** Deploy hybrid cryptographic schemes combining lightweight SMPC and TEE-based aggregation to encrypt and isolate update streams while maintaining acceptable latency. This limits interception and timing analysis by external or cross-institutional adversaries during synchronization over WAN links. **(4) Misconfigured audits/logging.** Integrate blockchain-based provenance tracking with metadata obfuscation and differentially private logging to ensure traceability without exposing participation frequency or institutional identifiers. This reduces internal leakage risks from logs and external deanonymization via audit metadata.

Table 4 provides the fine-tuning-phase taxonomy, capturing the federated update pathway, client/server adversaries, surface-specific threats, and the layered defenses discussed here.

## 5. Inference Phase

The vulnerabilities embedded during preprocessing and compounded through federated fine-tuning surface most clearly during inference, where PHI can leak through prompts, hidden states, caches, and outputs even when model parameters remain protected. In healthcare, LLM inferencing is the process where a trained LLM uses patient data to generate predictions, diagnoses, and personalized treatment recommendations in real time. However, this phase poses a distinct threat in LLM deployments, as clinical inputs can leak through outputs, embeddings, or memory without exposing model parameters. Protecting patient data during inference is essential for regulatory compliance and ethical AI.

### 5.1. Threat Model

Inference in healthcare LLM deployments introduces risks that differ from training. Here, PHI can leak via prompts, embeddings/hidden states, KV caches, and returned text, even when model weights and training data are protected. This matters for real clinical tools like discharge-note assistants, triage chatbots, radiology and pathology summarizers, and multilingual transcript systems—often deployed on hybrid edge–cloud stacks under tight latency constraints [2], [5], [6], [20], [27].

**5.1.1. Attacker Landscape and Incentives.** Inference introduces distinct adversarial incentives because sensitive information appears in runtime artifacts, *not* in training data or gradients. These artifacts persist in logs, caches, monitoring traces, and activation buffers maintained across large, distributed clinical infrastructures.

**Internal adversaries** pose the most persistent risk, as they interact directly with operational LLM systems deployed inside hospitals: **(1)** IT operators, SRE teams, and platform engineers often access telemetry dashboards, runtime logs, and GPU memory during debugging. Because multi-turn systems maintain *KV caches* across generations, these memory snapshots may reveal full conversations, including telemedicine notes or oncology assessments [24], [25]. **(2)** Clinicians and data scientists frequently input

full EHR excerpts (e.g., ICU timelines, staging notes, surgical histories) into prompts during validation or A/B testing. These prompts can resurface in logging, error traces, or monitoring pipelines [6], [30]. **(3)** Teams operating domain-specific inference pipelines, such as digital pathology or endocrine-cancer extraction systems, may store or inspect intermediate embeddings during performance tuning, inadvertently exposing PHI embedded in feature vectors [2]. Internal threats are especially severe because these activities occur under legitimate operational workflows and rarely trigger intrusion detection. External adversaries

**External adversaries** raise threats to exploit LLM deployment APIs, cross-hospital WAN links, or cloud inference infrastructure: **(1)** API-based attackers can perform *user inference attacks*, probing outputs to determine whether the model was adapted on rare clinical cohorts (e.g., rare oncology patients), enabling deanonymization of small hospitals [31], [32]. **(2)** Cloud-side or semi-trusted platform adversaries can apply *black-box inversion*, reconstructing sensitive spans from intermediate activations or returned outputs [88], [89]. **(3)** Network adversaries exploit WAN traffic between local hospitals and cloud inference nodes; timing and packet-size patterns can reveal PHI density or prompt structure, even when content is encrypted [90]. **(4)** Distributed inference architectures especially those using adapter offloading or split execution expand the attack surface when activations cross trust boundaries [33], [91]. These attackers operate outside institutional boundaries, motivated by financial, competitive, or strategic objectives.

**5.1.2. Prior Knowledge and Capability Gradient.** Adversaries at inference time vary in domain awareness and technical capability, shaping how effectively they exploit runtime artifacts such as prompts, KV caches, and hidden states. **(1)** Low-resource insiders, clinicians, analysts, or IT operators understand prompt formats and may access logs or traces containing raw PHI. Their local familiarity with EHR content and dashboard telemetry enables unintentional exposure through debugging or monitoring systems [25], [30]. **(2)** Moderate-capability actors, including engineers or cloud operators, know where GPU dumps, activation buffers, and inference logs are stored, allowing silent extraction or correlation of prompts, embeddings, and conversational states [24]. **(3)** High-resource adversaries such as external API probers or state-linked entities, possess auxiliary clinical datasets and use prior statistical knowledge to align output distributions or hidden-state signatures with real cohorts, revealing rare conditions or site-specific attributes [31], [32].

This knowledge gradient translates directly into exploitability: (1) insiders cause accidental PHI leakage via logs and KV caches; (2) mid-level actors conduct embedding inversion or output-correlation attacks using API or infrastructure access [29], [88]; and (3) high-resource adversaries perform cross-dataset alignment and timing analysis across WAN links to infer sensitive patterns or prompt structures [90], [91].

**5.1.3. Attack Surface Vulnerabilities and Enabled Attacks.** The key vulnerabilities include:

- **Prompt channel and Hidden-State Leakage:** Clinical prompts carry PHI such as diagnoses, medications, lab timelines, and family histories. Even when transmitted over TLS, prompts are plaintext at the service node, and simple redaction can distort clinical meaning. Hidden states store rich patterns from EHR notes or pathology text, which can be reconstructed by embedding inversion in semitrusted clouds or debug pipelines [30], [88], [89], [91]. Edge–cloud splits or adapter offloads can leak intermediate activations if not properly protected [33], [91].
- **KV cache leakage:** Multi-turn KV caches retain contextual information for efficiency. If snapshots of GPU memory or caches are exposed, entire patient dialogues, such as telemedicine or oncology sessions, can be recovered [24], [25]. Hybrid-cloud monitoring and scaling layers introduce additional points of exposure.
- **Returned text (outputs):** Outputs may leak PHI via over-specific recommendations (dose names, rare comorbidity patterns) or verbatim regurgitation of earlier prompts. Even with input obfuscation, outputs may re-expose sensitive facts if protections are not end-to-end [27], [29], [30].
- **Delegated / hybrid inference (edge–cloud):** When embeddings or adapters are computed locally and later processed in the cloud, activations and metadata traverse WAN links. Without secure partitioning or local DP, these streams can be profiled or linked to patient data [33], [91], [92]. WAN conditions can also create timing side channels that leak information [90].
- **Restoration/meta-vector channels:** Pipelines that remove sensitive spans and send noised restoration vectors risk leaking protected attributes (e.g., HIV status, hereditary cancer risk) if noise schedules are mismanaged [29]. Reusing obfuscation mappings across sessions increases this risk [30].
- **Operational logging and provenance:** LLM stacks with tracing, A/B testing, or auditing can inadvertently capture PHI in prompts, activations, or outputs. Without strict log minimization, this creates long-term leakage vulnerabilities [8], [25], [87].

These vulnerabilities enable concrete inference-time attack types highly relevant to healthcare: **(1) Prompt reconstruction attacks**, extracting PHI from hidden states, KV caches, or traced activations. **(2) Embedding inversion attacks**, reconstructing clinical notes, lab summaries, or pathology descriptions via inversion of hidden states or adapter outputs. **(3) Output-based inference attacks**, inferring patient cohort membership, rare disease participation, or underlying text from over-specific model responses. **(4) Split-inference intercept attacks**, recovering sensitive activations or adapter states crossing WAN links in hybrid edge–cloud deployments. **(5) Restoration-vector attacks**, recovering masked identifiers or sensitive attributes from poorly noised reconstruction vectors. **(6) Log-correlation attacks**, linking traces across clinical sessions to reconstruct longitudinal care histories. Together, these attacks operationalize the incentives and capabilities outlined above, turning routine inference pathways into high-value exploitation channels.

TABLE 5: Inference Phase — Summary of Threats, Defenses, Limitations, and Recommendations.

| Threat Model | Privacy-preserving Defenses | Limitations | Phase-Tied Recommendations |
|---|---|---|---|
| • **Internal adversaries:** IT/SRE teams with access to logs, GPU memory, KV caches; clinicians or data scientists entering full EHR prompts; teams inspecting intermediate embeddings.<br>• **External adversaries:** API-based attackers, cloud-side inversion agents, WAN interceptors, and adapter/split execution exploiters.<br>• **Capabilities:** prompt reconstruction, embedding inversion, output-based inference, cross-session correlation, KV-cache scraping, timing and metadata profiling.<br>• **Key Vulnerabilities:** prompts/inputs, hidden states, KV caches, returned text, edge–cloud activations, restoration/meta-vectors, operational logs/provenance traces. | • **Local DP:** DP-Forward, split-and-denoise, token perturbations (InferDPT).<br>• **Input Obfuscation/Span Removal:** PrivacyRestore, dynamic substitution/hashing, selective masking.<br>• **Cryptographic Protocols:** MPC-minimized, PermLLM (HE), GPU-accelerated FHE.<br>• **KV Cache Protection:** KV-Shield, PFID TEEs.<br>• **Model Partitioning/Edge–Cloud:** PFID local embeddings/adapters; PrivateLoRA for local adapter execution.<br>• **Federated Inference:** eFedLLM for distributed inference across institutions. | • **Prompt/Hidden-State Leakage:** LDP reduces fidelity; static obfuscation bypassed by cross-session correlation.<br>• **KV Cache Leakage:** relies on trusted hardware; persistent or misconfigured caches expose full dialogues.<br>• **Returned Text/Restoration:** span removal fails under repeated probing; restoration vectors reveal patient traits.<br>• **Hybrid Inference Leakage:** MPC/HE latency exposes timing channels; WAN activations profiled by network adversaries.<br>• **Operational Logging:** traces retain sensitive prompts/activations, creating persistent PHI leakage. | • **Prompt/Hidden-State Protection:** use span-level adaptive DP with contextual obfuscation for clinical entities.<br>• **Cross-Session Defense:** employ session-specific token randomization and dynamic obfuscation schedules.<br>• **KV Cache Security:** enforce automated expiration, TEE-based isolation, and strict memory hygiene.<br>• **Secure Hybrid Inference:** combine MPC with GPU-accelerated HE; run early layers locally before transmission.<br>• **Operational Logging/Provenance:** minimize trace capture; anonymize logs; use privacy-aware provenance tracking. |

## 5.2. Privacy-Preserving Defenses

Inference-time attacks differ from training risks by targeting live system outputs, embeddings, and caches - often assumed non-sensitive. In healthcare, this underscores the need for holistic, privacy-by-design defenses protecting both user inputs and transient data. These mechanisms aim to prevent leakage during runtime execution, when PHI is most exposed.

**Local Differential Privacy (LDP).** introduces calibrated noise to inputs or embeddings before model access, limiting adversarial inference. DP-Forward adds noise during the forward pass to obscure subtle prompt differences [89], while Split-and-Denoise combines token masking with local DP to resist reconstruction [93]. InferDPT adapts this for black-box inference with token-level perturbations [88]. For instance, a query like "blood in stool with recurring headaches" is perturbed locally, preserving ICD-prediction utility while masking identity cues. Excessive noise, however, may weaken rare-disease fidelity.

**Input Obfuscation and Privacy Span Removal.** selectively mask or substitute sensitive spans like names, dosages, or diagnoses before transmission. PrivacyRestore removes PHI locally, transmitting noised meta-vectors for server-side reconstruction [29]; Instance Obfuscation hashes tokens dynamically per session [30]. In clinical chatbots, such obfuscation ensures confidential prompts (e.g., "tested positive for HIV") never leave the device. When paired with local DP, these lightweight methods reinforce cross-session privacy while preserving reasoning quality, ideal for telemedicine and triage.

**Cryptographic Protocols (MPC and HE).** protect inference over untrusted clouds. MPC-minimized secret shares only early layers to hide embeddings efficiently [94]; PermLLM uses HE with lightweight permutations for secure attention under WAN latency [90]; and GPU-accelerated FHE supports encrypted inference for radiology or pathology batch tasks [1]. Though resource-intensive, such methods suit regulated domains like oncology where security outweighs delay.

**KV Cache Protection.** addresses persistent memory risks. KV-Shield permutes attention matrices to render stolen caches useless [24], while PFID confines cache-sensitive operations within Trusted Execution Environments (TEEs) [91]. These safeguards prevent full-session transcript recovery in telemedicine settings with minimal latency impact, assuming secure hardware is available.

**Model Partitioning and Edge-Cloud Collaboration.** localize sensitive processing by executing early layers or adapters on secure devices. PFID runs embeddings within TEEs before delegating deeper computation [91], while PrivateLoRA fine-tunes and executes low-rank adapters locally [33]. This allows mobile clinical apps to handle initial contexts (e.g., "47-year-old diabetic with chest pain") privately before cloud inference. While device demands increase, emerging mobile accelerators make shallow inference viable.

**Federated Inference and Decentralized Deployment.** extend FL principles to inference-time privacy. eFedLLM distributes inference across multiple nodes so no single server processes full prompts [92]. Used in oncology and multilingual telemedicine networks, each institution executes partial inference and contributes to aggregate predictions [4]. This decentralization reduces central memory risk but introduces synchronization challenges across heterogeneous hospital systems.

**Takeaway.** Inference-time threats arise not from raw data or gradients but from runtime artifacts, including prompts, hidden states, KV caches, intermediate activations, and returned text, where PHI may persist across sessions. Defenses such as Local DP, privacy-span removal, cryptographic inference (MPC/HE), KV-cache shielding, and edge–cloud partitioning align directly with the attack surfaces identified in the threat model. *Prompt and Hidden-State Leakage:* Local DP and input obfuscation prevent direct re-identification from plaintext prompts and hidden activations, reducing the risk of embedding in-

version or membership inference attacks. *KV-Cache Exposure:* TEEs and KV-Shield protect multi-turn dialogue caches stored in GPU memory, securing patient–clinician transcripts in telemedicine and ICU chat scenarios. *Output and Restoration Risks:* Privacy restore mechanisms remove or noise sensitive spans in returned text, limiting PHI resurfacing across prompts. *Delegated/Hybrid Inference:* Cryptographic protocols such as MPC and HE secure embeddings and activations traversing WAN links, shielding clinical data during cloud-based or split inference. Together, these defenses provide partial containment against prompt-level leakage, hidden-state reconstruction, and output correlation, enhancing privacy assurance across deployed healthcare LLMs.

**Limitation.** Despite their value, inference defenses remain fragile when mapped to real-world healthcare deployments: (1) *Prompt channel and Hidden-State Leakage)*: Local DP blurs sensitive fields but diminishes clinical fidelity; static obfuscation can still be bypassed via cross-session correlation. (2) *KV cache leakage*: KV shielding depends on trusted hardware; misconfigured or persistent caches allow recovery of multi-turn patient conversations. (3) *Returned text (outputs) and Restoration*: Span-level removal fails under repeated queries, allowing re-identification through prompt–output correlation. (4) *Delegated / hybrid inference*: MPC and HE secure content but introduce high latency, enabling timing-based inference of prompt complexity or PHI density. (5) *Operational logging and provenance*: Tracing and debugging pipelines often retain raw prompts and activations, creating durable PHI leakage paths even after session termination.

---

**Recommendation 3:** To address these shortcomings, inference-time privacy must adopt dynamic, context-aware protection tied to each identified surface:
**(1) Prompt and Hidden-State Protection**: Apply span-level adaptive DP and contextual obfuscation that selectively perturb only sensitive entities (e.g., HIV status, rare conditions), preserving clinical accuracy while blocking re-identification. **(2) Cross-session correlation Defense**: Introduce session-specific randomization of token mappings and obfuscation schedules to prevent linkage of repeated prompts across oncology or chronic-care cases. **(3) KV-Cache Security:** Enforce automated cache expiration and TEE-based isolation for multi-turn dialogue memory, preventing GPU scraping or reuse of prior session states. **(4) Secure Hybrid Inference:** Use hybrid cryptographic inference (MPC + GPU-accelerated HE) to balance latency and confidentiality across WAN-based deployments. **(5) Operational Logging and Provenance:** Minimize trace capture, anonymize diagnostic logs, and employ privacy-aware provenance tracking to limit long-term PHI persistence. Together, these measures form a phase-aware, layered strategy explicitly mapped to inference-time attack surfaces, closing the residual gaps between theoretical privacy and real-world healthcare deployment safety.

---

Table 5 presents the inference-phase taxonomy, detailing model-, prompt-, and system-level threats alongside their corresponding defensive strategies, limitations, and recommendations.

# 6. Cross-Phase Propagation of Privacy Risks

Although prior sections analyze vulnerabilities within preprocessing, federated fine-tuning, and inference individually, real-world healthcare deployments rarely operate in isolation. Privacy failures accumulate across phases, allowing seemingly minor leaks in one stage to amplify downstream. Preprocessing artifacts such as residual identifiers, demographic signals, or mislabeled data can be memorized during fine-tuning and later resurfaced through inference-time extraction attacks. Conversely, insecure inference interfaces (e.g., prompt injection, KV-cache probing) can exfiltrate model representations that encode sensitive patterns originating from earlier phases. This lifecycle coupling explains why phase-specific defenses: differential privacy in training, anonymization in preprocessing, or filtering in inference often fail when deployed independently.

Propagation also creates indirect vulnerabilities. Poisoned or synthetic records introduced during preprocessing can bias model gradients during fine-tuning, which then distort inference behavior in clinically consequential ways. Similarly, unsecured logs, caching mechanisms, or checkpoint reuse allow auxiliary metadata from later stages to re-identify samples that were previously anonymized. Threats therefore move bidirectionally: upstream preprocessing influences representational leakage downstream, while downstream inference interfaces expose training-time weaknesses upstream. This interconnectedness underscores a key observation of this SoK: protecting any single phase without accounting for its interactions with the others is structurally insufficient.

This insight motivates the need for lifecycle-aware privacy mechanisms, which the conclusion elaborates by integrating cross-phase findings into a coherent set of deployment-ready recommendations.

# 7. Conclusion

This work presents the first phase-aware systematization of privacy and security in healthcare LLMs, tracing risks and defenses across data preprocessing, fine-tuning, and inference. By mapping adversaries, capabilities, attack surfaces, and defenses at each stage, we illustrate how privacy risks manifest uniquely across the LLM lifecycle and why phase-agnostic approaches fall short in clinical settings

Our analysis highlights a central gap: existing techniques, from anonymization and DP to secure aggregation and cryptographic inference, provide protection only in isolation. In practice, they struggle with **rare-disease exposure**, **multi-round reconstruction**, **poisoning**, and **runtime leakage** through logs, caches, and distributed infrastructure. Effective privacy in healthcare LLMs therefore requires defenses that are **context-aware**, **workflow-aligned**, and **robust across phases**, not just individually strong. We identify three priorities for future research: (1) **layered, adaptive defenses** that combine DP, randomized adapters, and hardware enclaves with phase-specific tuning; (2) **standardized evaluation and auditing protocols** for jointly assessing privacy and clinical fidelity; and (3)

**lightweight secure computation and federated inference** that meet the latency and reliability demands of real clinical deployments.

By structuring a fragmented landscape into a unified, phase-aware framework, this SoK provides a clear roadmap for developing privacy-resilient, regulation-aligned, clinically safe LLM systems that can earn trust in high-stakes healthcare settings.

# References

[1] L. de Castro, A. Polychroniadou, and D. Escudero, "Privacy-preserving large language model inference via gpu-accelerated fully homomorphic encryption," *Neurips Safe Generative AI Workshop 2024*, 2024.

[2] D. T. Lee, A. Vaid, K. M. Menon, R. Freeman, D. S. Matteson, M. P. Marin, and G. N. Nadkarni, "Development of a privacy preserving large language model for automated data extraction from thyroid cancer pathology reports," 2023, medRxiv, 2023-11.

[3] I. C. Wiest, D. Ferber, J. Zhu, M. van Treeck, S. K. Meyer, R. Juglan, and J. N. Kather, "Privacy-preserving large language models for structured medical information retrieval," *NPJ Digital Medicine*, vol. 7, no. 1, p. 257, 2024.

[4] R. Ye, W. Wang, J. Chai, D. Li, Z. Li, Y. Xu, Y. Du, Y. Wang, and S. Chen, "Openfedllm: Training large language models on decentralized private data via federated learning," *Association for Computing Machinery*, 2024.

[5] A. Manoel, M. Garcia, T. Baumel, S. Su, J. Chen, R. Sim, D. Miller, D. Karmon, and D. Dimitriadis, "Federated multilingual models for medical transcript analysis," *Conference on Health, Inference, and Learning, 2023*, 2023.

[6] Z. Wang, H. Li, D. Huang, and A. M. Rahmani, "Healthq: Unveiling questioning capabilities of llm chains in healthcare conversations," 2024, arXiv preprint arXiv:2409.19487.

[7] N. Khalid, A. Qayyum, M. Bilal, A. Al-Fuqaha, and J. Qadir, "Privacy-preserving artificial intelligence in healthcare: Techniques, challenges and future research directions," *Computers in Biology and Medicine*, 2023.

[8] J. Jonnagaddala and Z. S. Y. Wong, "Privacy preserving strategies for electronic health records in the era of large language models," *npj Digital Medicine*, vol. 8, no. 1, p. 34, 2025.

[9] A. K. Conduah, S. Ofoe, and D. Siaw-Marfo, "Data privacy in healthcare: Global challenges and solutions," *Frontiers in Digital Health*, 2025.

[10] O. Shoham and N. Rappoport, "Federated learning of medical concepts embedding using behrt," *JAMIA Open. 2024*, 2024.

[11] M. Aljohani, J. Hou, S. Kommu, and X. Wang, "A comprehensive survey on the trustworthiness of large language models in healthcare," 2025.

[12] M. S and S. MJ., "Large language models in healthcare and medical applications: A review," *Bioengineering*, 2025.

[13] *Supervised learning.* Elsevier eBooks, 2022.

[14] J. Thomas, *Preprocessing.* Informa, 2023, pp. 196–210.

[15] D. Liu and T. Miller, "Federated pretraining and fine tuning of bert using clinical notes from multiple silos," 2020, arXiv preprint arXiv:2002.08562.

[16] C. Christophe, P. K. Kanithi, P. Munjal, T. Raha, N. Hayat, R. Rajan, A. Al-Mahrooqi, A. Gupta, M. U. Salman, G. Gosal *et al.*, "Med42–evaluating fine-tuning strategies for medical llms: full-parameter vs. parameter-efficient approaches," *arXiv preprint arXiv:2404.14779*, 2024.

[17] F. Lucini, "The real deal about synthetic data," *MIT Sloan Management Review*, vol. 63, no. 1, pp. 1–4, 2021.

[18] H. Yang, H. Chen, H. Guo, Y. Chen, C. S. Lin, S. Hu, and X. Wang, "Llm-medqa: Enhancing medical question answering through case studies in large language models," 2024, arXiv preprint arXiv:2501.05464.

[19] H. Lu and U. Naseem, "Can large language models enhance predictions of disease progression? investigating through disease network link prediction," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 17703–17715.

[20] H. Jung, Y. Kim, H. Choi, H. Seo, M. Kim, J. Han, and Y. H. Kim, "Enhancing clinical efficiency through llm: Discharge note generation for cardiac patients," 2024, arXiv preprint arXiv:2404.05144.

[21] Gagan, N, Sanand, and Sasidharan, "Enhancing oncology care with federated learning and foundation models," *2024 ITU Kaleidoscope: Innovation and Digital Transformation for a Sustainable World (ITU K)*, 2024.

[22] Y. Li, Z. Li, K. Zhang, R. Dan, S. Jiang, and Y. Zhang, "A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge." *Cureus. 2023*, 2023.

[23] A. Anaissi, A. Braytee, and J. Akram, "Fine-tuning llms for reliable medical question-answering services," 2024, arXiv preprint arXiv:2410.16088.

[24] H. Yang, D. Zhang, Y. Zhao, Y. Li, and Y. Liu, "A first look at efficient and secure on-device llm inference against kv leakage," *Proceedings of the 19th Workshop on Mobility in the Evolving Internet Architecture, 2024*, 2024.

[25] D. Chen, A. Youssef, R. Pendse, A. Schleife, B. K. Clark, H. Hamann, J. He *et al.*, "Transforming the hybrid cloud for emerging ai workloads," *arXiv*, 2025.

[26] D. Oniani, X. Wu, S. Visweswaran, S. Kapoor, S. Kooragayalu, K. Polanska, and Y. Wang, "Enhancing large language models for clinical decision support by incorporating clinical practice guidelines," in *2024 IEEE 12th International Conference on Healthcare Informatics (ICHI)*. IEEE, 2024, pp. 694–702.

[27] A. Ghosh, A. Acharya, R. Jain, S. Saha, A. Chadha, and S. Sinha, "Clipsyntel: clip and llm synergy for multimodal question summarization in healthcare," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 20, 2024, pp. 22031–22039.

[28] Z. Wang, Z. Shen, Y. He, G. Sun, H. Wang, L. Lyu, and A. Li, "Flora: Federated fine-tuning large language models with heterogeneous low-rank adaptations," *arXiv preprint arXiv:2409.05976*, 2024.

[29] Z. Zeng, J. Wang, J. Yang, Z. Lu, H. Zhuang, and C. Chen, "Privacyrestore: Privacy-preserving inference in large language models via privacy removal and restoration," *arXiv preprint arXiv:2406.01394, 2024*, 2024.

[30] Y. Yao, F. Wang, S. Ravi, and M. Chen, "Privacy-preserving language model inference with instance obfuscation," *arXiv preprint arXiv:2402.08227, 2024*, 2024.

[31] N. Kandpal, K. Pillutla, A. Oprea, P. Kairouz, C. A. Choquette-Choo, and Z. Xu, "User inference attacks on large language models," *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024.

[32] R. Staab, M. Vero, M. Balunović, and M. Vechev, "Beyond memorization: Violating privacy via inference with large language," *arXiv preprint arXiv:2310.07298, 2023*, 2023.

[33] Y. Wang, Y. Lin, X. Zeng, and G. Zhang, "Privatelora for efficient privacy preserving llm," *arXiv preprint arXiv:2311.14030, 2023*, 2023.

[34] A. Gadotti, L. Rocher, F. Houssiau, A. M. Crețu, and Y. A. D. Montjoye, "Anonymization: The imperfect science of using data while preserving privacy," *Science Advances*, vol. 10, no. 29, p. eadn7053, 2024.

[35] R. Tertulino, N. Antunes, and H. Morais, "Privacy in electronic health records: a systematic mapping study," *Journal of Public Health*, vol. 32, no. 3, pp. 435–454, 2024.

[36] M. Hussain, N. Akhtar, and R. Hasan, "A robust framework for ensuring data confidentiality and security in modern healthcare networks," *International Journal of Intelligent Systems and Applications in Engineering*, vol. 12, no. 3, pp. 128–137, 2024.

[37] R. Zhang, R. Xue, and L. Liu, "Searchable encryption for healthcare clouds: A survey," *IEEE Transactions on Services Computing*, vol. 11, no. 6, pp. 978–996, 2017.

[38] S. Abdali, R. Anarfi, C. J. Barberan, and J. He, "Securing large language models: Threats, vulnerabilities and responsible practices," *arXiv preprint arXiv:2403.12503*, 2024.

[39] Anonymous, "Ensuring data security and compliance in etl processes for healthcare and financial services," *ResearchGate (preprint)*, 2024.

[40] J. Cândido, M. Aniche, and A. V. Deursen, "Log-based software monitoring: a systematic mapping study," *PeerJ Computer Science*, vol. 7, p. e489, 2021.

[41] H. Bahsi and A. Levi, "Preserving organizational privacy in intrusion detection log sharing," in *2011 3rd International Conference on Cyber Conflict*. IEEE, June 2011, pp. 1–14.

[42] National Audit Office, "Investigation: Wannacry cyber attack and the nhs," National Audit Office, UK, Tech. Rep., 2018.

[43] F. Wang and B. Li, "Data reconstruction and protection in federated learning for fine-tuning large language models," *IEEE Transactions on Big Data, 2024*, 2024.

[44] H. Li, Y. Chen, J. Luo, J. Wang, H. Peng, Y. Kang, and Y. Song, "Privacy in large language models: Attacks, defenses and future directions," 2023, arXiv preprint arXiv:2310.10383.

[45] E. Shayegani, M. A. A. Mamun, Y. Fu, P. Zaree, Y. Dong, and N. Abu-Ghazaleh, "Survey of vulnerabilities in large language models revealed by adversarial attacks," *arXiv preprint arXiv:2310.10844*, 2023.

[46] A. Narayanan and V. Shmatikov, "Robust de-anonymization of large sparse datasets," in *2008 IEEE Symposium on Security and Privacy (sp 2008)*. IEEE, 2008, pp. 111–125.

[47] T. Stadler, B. Oprisanu, and C. Troncoso, "Synthetic data–anonymisation groundhog day," in *31st USENIX Security Symposium (USENIX Security 22)*. USENIX Association, 2022, pp. 1451–1468.

[48] "Cyber threat modeling of an llm-based healthcare system," *International Conference on Information Systems Security and Privacy*, 2025.

[49] M. Jagielski, A. Oprea, B. Biggio, C. Liu, C. Nita-Rotaru, and B. Li, "Manipulating machine learning: Poisoning attacks and countermeasures," in *Proceedings of the IEEE Symposium on Security and Privacy (SP)*, 2018, pp. 19–35.

[50] W. F. Shah, "Data preprocessing in healthcare: A vital step towards informed decision-making," 2024.

[51] E. Shamsinejad, T. Banirostam, M. M. Pedram, and A. M. Rahmani, "A review of anonymization algorithms and methods in big data," *Annals of Data Science*, pp. 1–27, 2024.

[52] O. Vovk, G. Piho, and P. Ross, "Methods and tools for healthcare data anonymization: a literature review," *International Journal of General Systems*, vol. 52, no. 3, pp. 326–342, 2023.

[53] A. Aminifar, Y. Lamo, K. I. Pun, and F. Rabbi, "A practical methodology for anonymization of structured health data," 2019.

[54] S. L. Ribeiro and E. T. Nakamura, "Privacy protection with pseudonymization and anonymization in a health iot system: results from ocariot," in *2019 IEEE 19th International Conference on Bioinformatics and Bioengineering (BIBE)*. IEEE, 2019, pp. 904–908.

[55] I. E. Olatunji, J. Rauch, M. Katzensteiner, and M. Khosla, "A review of anonymization for healthcare data," *Big data*, vol. 12, no. 6, pp. 538–555, 2024.

[56] Y. Qu, Y. Dai, S. Yu, P. Tanikella, T. Schrank, T. Hackman, and D. Wu, "A novel compact llm framework for local, high-privacy ehr data applications," 2024, arXiv preprint arXiv:2412.02868.

[57] M. Fredrikson, S. Jha, and T. Ristenpart, "Model inversion attacks that exploit confidence information and basic countermeasures," in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 2015, pp. 1322–1333.

[58] D. R. Alattal, Z. Wang, P. Myles, and A. Tucker, "Creating synthetic geospatial patient data to mimic real data whilst preserving privacy," in *2023 IEEE 36th International Symposium on Computer-Based Medical Systems (CBMS)*. IEEE, 2023, pp. 7–12.

[59] M. Jin, Q. Yu, C. Zhang, D. Shu, S. Zhu, M. Du, and Y. Meng, "Health-llm: Personalized retrieval-augmented disease prediction model," 2024, arXiv preprint arXiv:2402.00746.

[60] J. A. Alzate-Grisales, J. Bernal-Salcedo, J. M. Saborit-Torres, A. Mora-Rubio, J. M. Serrano, F. García-García, and M. D. L. Iglesia-Vayá, "Dismed-llm: De-identifying spanish medical text with large language models," 2025.

[61] M. Abbasian, I. Azimi, A. M. Rahmani, and R. Jain, "Conversational health agents: A personalized llm-powered agent framework," 2023, arXiv preprint arXiv:2310.02374.

[62] M. A. Roshani, X. Zhou, Y. Qiang, S. Suresh, S. Hicks, U. Sethuraman, and D. Zhu, "Generative llm powered conversational ai application for personalized risk assessment: A case study in covid-19," 2024, arXiv preprint arXiv:2409.15027.

[63] M. Naji, M. Masmoudi, and H. B. Zghal, "Towards an llm based approach for medical e-consent," *Procedia Computer Science*, vol. 246, pp. 3694–3701, 2024.

[64] A. Goncalves, P. Ray, B. Soper, J. Stevens, L. Coyle, and A. P. Sales, "Generation and evaluation of synthetic patient data," *BMC medical research methodology*, vol. 20, pp. 1–40, 2020.

[65] X. Chen, Z. Wu, X. Shi, H. Cho, and B. Mukherjee, "Generating synthetic electronic health record (ehr) data: A review with benchmarking," 2024, arXiv preprint arXiv:2411.04281.

[66] I. Pérez-Díez, R. Pérez-Moraga, A. López-Cerdán, J. M. Salinas-Serrano, and M. D. la Iglesia-Vayá, "De-identifying spanish medical texts-named entity recognition applied to radiology reports," *Journal of Biomedical Semantics*, vol. 12, pp. 1–13, 2021.

[67] M. Ibrahim, Y. A. Khalil, S. Amirrajab, C. Sun, M. Breeuwer, J. Pluim, and M. Dumontier, "Generative ai for synthetic data across multiple medical modalities: A systematic review of recent developments and challenges," 2024, arXiv preprint arXiv:2407.00116.

[68] M. Ali, M. Ali, M. Hussain, and D. Koundal, "Generative adversarial networks (gans) for medical image processing: Recent advancements," *Archives of Computational Methods in Engineering*, pp. 1–14, 2024.

[69] S. Kweon, J. Kim, J. Kim, S. Im, E. Cho, S. Bae, and E. Choi, "Publicly shareable clinical large language model built on synthetic clinical notes," 2023, arXiv preprint arXiv:2309.00237.

[70] Z. Wang, P. Myles, and A. Tucker, "Generating and evaluating cross-sectional synthetic electronic healthcare data: preserving data utility and patient privacy," *Computational Intelligence*, vol. 37, no. 2, pp. 819–851, 2021.

[71] K. M. Babu, E. Bhavitha, M. Mythri, A. Anusha, B. S. Chandana, and C. G. Akhtar, "Privacy-preserving federated learning for healthcare: A synergistic approach using differential privacy and homomorphic encryption," 2024, available at SSRN 5088943.

[72] S. Sadeepa, K. Kavinda, E. Hashika, C. Sandeepa, T. Gamage, and M. Liyanage, "Disllm: Distributed llms for privacy assurance in resource-constrained environments," in *2024 IEEE Conference on Communications and Network Security (CNS)*. IEEE, 2024, pp. 1–9.

[73] T. Bossy, J. Vignoud, T. Rabbani, J. R. T. Pastoriza, and M. Jaggi, "Mitigating unintended memorization with lora in federated learning for llms," 2025, arXiv preprint arXiv:2502.05087.

[74] Kokala and Abhilash, "Scalable large language models for the healthcare domain: A research perspective," 2025.

[75] H. Fu, H. Chen, S. Lin, and G. Xing, "Shade-ad: An llm-based framework for synthesizing activity data of alzheimer's patients," 2025, arXiv preprint arXiv:2503.01768.

[76] F. Piccialli, D. Chiaro, P. Qi, V. Bellandi, and E. Damiani, "Federated and edge learning for large language models." *Information Fusion*, 2025.

[77] Liu, Xiao-Yang, Zhu, Rongyi, Zha, Daochen, Gao, Jiechao, Zhong, Shan, White, Matt, Qiu, and Meikang, "Differentially private low-rank adaptation of large language model using federated learning," *Association for Computing Machinery*, 2024.

[78] C. Li, B. Gu, Z. Zhao, Y. Qu, G. Xin, J. Huo, and L. Gao, "Federated transfer learning for on-device llms efficient fine tuning optimization," *Big Data Mining and Analytics*, 2025.

[79] X. Zuo, M. Wang, T. Zhu, S. Yu, and W. Zhou, "Large language model federated learning with blockchain and unlearning for cross-organizational collaboration," *arXiv preprint*, 2024.

[80] Z. Lin, X. Hu, Y. Zhang, Z. Chen, Z. Fang, X. Chen, A. Li, P. Vepakomma, and Y. Gao, "Splitlora: A split parameter-efficient fine-tuning framework for large language models," *arXiv preprint*, 2024.

[81] J. Zhao, "Privacy-preserving fine-tuning of artificial intelligence (ai) foundation models with federated learning, differential privacy, offsite tuning, and parameter-efficient fine-tuning (peft)," *Authorea Preprints, 2023*, 2023.

[82] B. C. Das, M. H. Amini, and Y. Wu, "Security and privacy challenges of large language models: A survey," *ACM Computing Surveys*, vol. 57, no. 6, pp. 1–39, 2025.

[83] J. Jiang, H. Jiang, Y. Ma, X. Liu, and C. Fan, "Low-parameter federated learning with large language models," *International Conference on Web Information Systems and Applications*, 2024.

[84] F. Jiang, L. Dong, S. Tu, Y. Peng, K. Wang, K. Yang, C. Pan, and D. Niyato, "Personalized wireless federated learning for large language models," *arXiv preprint*, 2024.

[85] F. Wu, Z. Li, Y. Li, B. Ding, and J. Gao, "Fedbiot: Llm local fine-tuning in federated learning without full model," *Association for Computing Machinery*, 2024.

[86] J. Qi, Z. Luan, S. Huang, C. Fung, H. Yang, and D. Qian, "Fdlora: Personalized federated learning of large language model via dual lora tuning," *arXiv preprint arXiv:2406.07925*, 2024.

[87] Z. Fang, Z. Lin, Z. Chen, X. Chen, Y. Gao, and Y. Fang, "Automated federated pipeline for parameter-efficient fine-tuning of large language models," *arXiv preprint*, 2024.

[88] M. Tong, K. Chen, J. Zhang, Y. Qi, W. Zhang, N. Yu, T. Zhang, and Z. Zhang, "Inferdpt: Privacy-preserving inference for black-box large language models," *IEEE Transactions on Dependable and Secure Computing, 2025*, 2025.

[89] M. Du, X. Yue, S. Chow, T. Wang, C. Huang, and H. Sun, "Dp-forward: Fine-tuning and inference on language models with differential privacy in forward pass," *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications, 2023*, 2023.

[90] F. Zheng, C. Chen, Z. Han, and X. Zheng, "Permllm: Private inference of large language models within 3 seconds under wan," *arXiv preprint arXiv:2405.18744, 2024*, 2024.

[91] H. Yang, Z. Li, Y. Zhang, J. Wang, N. Cheng, M. Li, and J. Xiao, "Pfid: Privacy first inference delegation framework for llms," *arXiv preprint arXiv:2406.12238, 2024*, 2024.

[92] S. Ding and C. Hu, "efedllm: Efficient llm inference based on federated learning," *arXiv preprint arXiv:2411.16003, 2024*, 2024.

[93] P. Mai, R. Yan, Z. Huang, Y. Yang, and Y. Pang, "Split-and-denoise: Protect large language model inference with local differential privacy," *arXiv preprint arXiv:2310.09130, 2023*, 2023.

[94] D. Rathee, D. Li, I. Stoica, H. Zhang, and R. Popa, "Mpc-minimized secure llm inference," *arXiv preprint arXiv:2408.03561, 2024*, 2024.