(RESEARCH ARTICLE)

# Trust, but Verify: Audit-ready logging for clinical AI

Jimmy Joseph *

*Solutions Engineer Advisor Sr., United states.*

## Abstract

High-end clinical AI systems today (e.g. in radiology and oncology) need to be transparent so that their final output (or recommendations) can be explained post hoc in detail. Workflows for audit-capable logging in medical imaging AI are available online [4] (IHE BALP). Highlights - We meet the clinical need for auditability with tamper evident logging integrated into everyday PACS and EHR workflows. We algorithmically timestamp, through cryptographically auditable operations, the entirety of data flow and model inferences in our approach using standard healthcare data formats (DICOM and HL7/FHIR), binding events in hash chains and Merkle trees. We assess the system using real and synthetic datasets, as well as simulated audit trails of tampering. The logging pipeline immediately performs runtime canonicalization of events, hashing, and periodic notarization to an external trusted service for anchoring immutable integrity. In experiments, our system detected ≈100% of unauthorized log modifications (95% CI reported in Results) and achieved an end-to-end latency of <5 ms for detecting tampering attempts (modifications, reordering, deletions, forgeries). The median time to verify evidence of integrity breaches was 3.2 ms per event, with no throughput impact (well over 10,000 events/s ingestion). Statistical analysis demonstrates excellent performance: sensitivity ≈100%, specificity ≈100%, and F1 ≈0.999 (95% CI in Results) for tamper detection, achieving reliable results across diverse hospital sites and device vendors. Storage overhead was <5% of total data size (≈120 bytes per log event). These findings show that our audit-ready logging can offer strong forensic trails for clinical AI, ensuring reproducibility and compliance (HIPAA, FDA 21 CFR Part 11) without causing delays in clinical workflows. We contend that such verifiable logs will be significant for clinical incident investigation, liability reduction, and meeting new regulatory constraints aimed at AI traceability (e.g. EU AI Act Article 12). By rendering AI decisions completely auditable and provably immutable, this work establishes a basis for safer and more responsible deployment of AI in patient care.

**Keywords:** Clinical Artificial Intelligence; Audit-Ready Loggin; Tamper-Evident Systems; Provenance and Reproducibility; Regulatory Compliance (HIPAA, EU AI Act), Forensic Readiness in Healthcare IT

## 1. Introduction

In high-stakes clinical domains such as radiology and oncology, as artificial intelligence (AI) becomes more commonly deployed, confidence in those systems has become very important. (see Fehr et al. [1]). Here, trust incorporates the capability to trace and verify the origin of model outputs, audit decision paths, and replicate results under scrutiny. In actual deployment, this means effectively keeping the trail of provenance that allows one to trace and understand the data and the predictions made on it - from the imaging modality and PACS, via the inference of the AI model, to the system such as an EHR or FHIR-based system where the result is registered. This level of auditability is not only important for clinician confidence but is also fast becoming required by regulators and governance regimes for safe AI. As an example, HIPAA already requires strong access logging (audit trails) and long (≥6 years) data retention in order to support breach investigation. Similarly, the draft EU AI Act [9] mandates that high-risk AI systems are traceable with logs of what they do automatically, and such logs must be resistant to tampering. In this environment, an audit-ready

* Corresponding author: Jimmy Joseph.

logging system acts as an important safety net: it produces a clear record of how an AI arrived at a recommendation and indicates if anything unusual happened in that process.

Yet, while this is the case, in many clinical AI deployments there are few good audit trails beyond simple transaction logs. In a standard radiology workflow, an imaging study is performed and saved into a PACS system, an AI algorithm automatically processes a set of images (e.g., a nodule is detected or a lesion of concern is marked) and delivers the result to the radiologist's worklist or directly injects it into a report in an EHR. Figure 1 depicts a PACS ->AI -> EHR integration involving three systems communicating with each other. Although protocols such as DICOM, HL7, and FHIR facilitate exchange between these components, they do not necessarily ensure auditability and integrity of the workflow. There are a number of ways that failures can and do manifest at various stages - a model drift over time may result in poor performance with new patient populations, data leakage or mishandling might put patient privacy at risk, or adversaries may feed the model adversarial inputs (e.g., slightly low-resolution images) to cause the model to fail. There exist demonstrations of adversarial attacks resulting in diagnostic AI misclassifying medical images with no apparent hint. Incidents like these are why we need "forensic" capabilities - such as detailed logs to interrogate what the model saw and did, whether data distribution shift was detected, and whether any anomaly (natural or malicious) occurred at processing time.

Audit logs are a well-recognized feature in clinical IT used for security and compliance. For example, the IHE ATNA (Audit Trail and Node Authentication) profile and DICOM Part 15 standard specify how healthcare systems should log security-sensitive events (such as user logins, data exchanges, etc.). These logs generally keep track of who looked at which patient data, when, and what action was performed. However, the existing audit trails handle user and system access events and do not reveal how an AI algorithm made decisions inside. To address this, in a clinical AI setting we would need to extend auditing to model inference events (e.g., associating the generated label with a single image or dataset, version and parameters of the model, any post-processing applied, or decision context enabled by a human in the loop). We also need it to be tamper-proof: once it is written, any user (including system administrators) should not be able to change or delete an event without someone noticing. This is essential for "trust but verify" - if clinicians and regulators trust the system to work, they need the system to work both correctly and honestly to verify it in retrospect.

In this paper, we propose a fully audit-ready logging framework for clinical AI systems, addressing the provenance, integrity, and reproducibility of model outputs. Our contributions are as follows:

Cryptographically Verifiable Logging. We develop a logging pipeline that records each stage of the clinical AI workflow (data ingestion, model inference, and output reporting) and secures it cryptographically, with hash linking (hash chains and Merkle trees), to allow detection of any tampering or omission with near-zero false positives. Log records consist of core metadata (patient & study IDs, model IDs, timestamps) in a canonical format to ensure that hashes agree between systems.

Interoperability Standards-Based. The architecture uses healthcare de facto standards (FHIR AuditEvent [2], DICOM audit messages, HL7 v2) for transparent integration within clinical systems. We propose a canonical JSON schema for audit events to capture DICOM image pointer data, HL7 messages, and FHIR resources, thus promoting interoperability while posing minimal encumbrance to legacy PACS/EHR environments. It also enables connecting our logs to hospital Security Information and Event Management (SIEM) systems and IHE ATNA repositories.

Tamper-Evidence Using Hybrid Merkle Chain. We propose a hybrid form of integrity with linear hash chaining and periodic Merkle tree checkpoints notarized by an external party. This hybrid construction yields efficient log segment proofs in $O(\log n)$ and strong security (even if a privileged attacker completely subverts the logging server, history cannot be rewritten without invalidating a cryptographic integrity predicate). We give a formal definition of the predicate and bound false negative tamper detection.

Evaluation on Simulated Clinical Settings. We synthesize an audit log dataset to simulate the log of a multi-site radiology deployment, including "normal" operational behaviors and injected tampering activities (e.g., tampered results, log deletion, timestamp alterations, etc.). We compare the detection rates (sensitivity/specificity), performance overhead (latency per log entry, throughput), and storage overhead. Here are some headline results: ~99.9% tamper detection and near-zero false alarm rate, with verification times of less than 1 millisecond in many cases. We show forensic queries, such as reconstructing image hashes of an input to an AI model from the log so that the original input can be checked for correctness.

Governance and Compliance Artifacts. Beyond technical findings, we provide practical artifacts for responsible AI governance: a Bias & Generalizability plan for the logging coverage, a Model Card for AI model descriptions and intended

use, an Audit Card identifying its scope and compliance (similar to a security "nutrition label"), a rigorous Security Threat Model considering adversarial threat models, and an AHIMA Privacy/Compliance checklist aligned with HIPAA and GDPR compliance requirements. These materials are intended to assist practitioners and regulators in implementing our approach in practical applications.

In summary, our work fills a vital gap between the potential of AI in health and the practical protections required for its safe use. By giving clinical AI a "black box recorder" that renders it audit-ready, we enable trust through transparency: we can trust what the AI says because we can subsequently check what went on inside the black box. Next, we describe related work in medical audit logging, ML provenance, and tamper-evident systems. We then explain our method and system design, present an implementation and comprehensive evaluation, and conclude with a discussion of the influence on clinical forensics, security, and future regulatory integration.

## 2. Related Work

### 2.1. Audit Logging in Healthcare IT

Audit logs for access to and usage of sensitive data have been a standard practice for health information systems for security and compliance reasons. One of the profiles in the IHE IT Infrastructure is the IHE ATNA [4] (Audit Trail and Node Authentication), which defines a standard to govern secure transmission and storage of audit messages between clinical systems. ATNA uses a DICOM-defined audit message format (see DICOM PS3.15, Annex A.5) and commonly leverages syslog over TLS to forward audit records to a centralized Audit Record Repository. The DICOM [3] standard includes an audit trail profile with a common schema for events, such as ActiveParticipant (user or process), EventID/Action (what was done), Patient/Object info, etc., ensuring that logs contain the essential information: 'who,' 'which,' 'when,' 'from where,' and 'doing what' activity (create, read, update, delete, etc.), who did it, and how. Such audit trails are an extremely important feature for privacy monitoring and accountability. They can be used, for example, to identify inappropriate access to the EHR system or to examine what was done after an incident.

### 2.2. Provenance and Reproducibility in ML Models

In non-healthcare fields, a standard for data provenance - the record of both the origin and subsequent processing history of data - has been developed and would be directly applicable to AI audit trails. The W3C PROV model defines an ontology to represent entities (data items), activities (processes), agents, and their relationships (wasDerivedFrom, wasGeneratedBy, etc.). For example, PROV can record that Report_X was produced by AI_Algorithm_Y using Image_Z. Provenance tracking is crucial for reproducibility in data science: if we know exactly what data and code contributed to a conclusion, we have a better chance of reproducing or auditing that conclusion. In the lifecycle of ML, provenance metadata typically includes dataset versions, model hyperparameters, software environments, and dependency lineage. Software solutions such as MLflow, Pachyderm, and DVC (Data Version Control) generate log entries for experiments and model artifacts with the ability to trace those entities by unique IDs and hashes. To track machine learning workflows, researchers have suggested extensions of the PROV model (e.g., ProvML or MLProv) to represent data preprocessing, model training, or inference, not as loosely structured information. For instance, DLProv is a provenance capture toolkit for deep learning that records all performed steps in W3C PROV and automatically reports a full trace of the entire model development process [12].

Provenance is equally important in the context of clinical AI: we want to know which model (including version and underlying training data) produced a particular output, because models can (and do) get updated or retrained from time to time. If an AI gets a case wrong, provenance data could provide insight into whether it relied on an outdated model or whether there was a problem in the data pipeline. Reproducibility is also related to provenance - can you rerun the analysis on the same data with the same model and get the same result? Regulators such as the FDA have called out reproducibility and transparency for AI/ML as components of Good Machine Learning Practice (GMLP). [13] For example, one FDA guideline requires tracking of the training and testing environments for the model to enable subsequent auditing of its performance on various subpopulations. Audit logs in our system include reminders of ML provenance: each log entry of an inference can point to the model's unique ID/hash, to a version of the training dataset, to validation performance metrics (by linking to the Model Card), and to information about applied preprocessing. This enables retrospective analysis to answer questions such as "Was the AI employing the latest approved model? What data was it trained on? Did we drift from training statistics?" - all crucial in building trust and accountability.

One problem is that provenance metadata can be quite large. We avoid this by not placing full models or datasets in the log but instead hashes or IDs of them (e.g., model hash, dataset ID in a repository). This follows the content-addressable storage idea in data lineage: you can refer to a large object by its short hash and fetch and verify the object later.

Concretely, we build on provenance notions from MLOps to provide an audit log that connects data -> model -> result. This surpasses plain IT logging and enters into algorithmic transparency as an important component of responsible AI. The key difference is that we convert these provenance records into an unalterable log trail, not just an experiment tracking database, with a security flavor as discussed next.

## 2.3. Tamper-Evident Logging and Integrity Techniques

The idea of tamper-evident logs goes back to computer security research in the 1990s. The desired result is logging systems that, even when compromised, cannot have their history altered without generating cryptographic "breadcrumbs." One basic scheme is based on hash chains (also referred to as linked hashes or hash lists). In its simplest form, a hash chain places the hash of the previous entry into the current one; taken to the extreme, it's the hash of the hash of the hash of the previous entry… all the way down to the first event. Should an attacker tamper with (or delete) an entry in the middle, the chain is no longer continuous - a future verifiable hash will not match. Schneier and Kelsey [6] proposed secure audit logs as a combination of not only hash chains and frequent digital signatures but also forward-secure encryption of log records so that an attacker who compromises the system in the future cannot read past log information. In their model, the logger calculates $h_i = H!\left(h_{i-1} \mid \mathrm{canon}(E_i)\right)$ where each new event is added to the accumulator, and the logger logs a signed checksum of the log after each update. This way, an attacker who takes control of the logger cannot retrospectively change what events occurred without failing to recreate the previous signature (as they would not have had the signing key at that time). This concept inspired a number of works, such as secure logging appliances and RFC proposals for signed syslog.

One drawback of naive hash chains is that it can take significant time to check for inclusion of a single event or a set of events because it may be necessary to scan the entire log. For instance, proving that something that happened a year ago is intact may require hashing over millions of entries. Merkle trees solve this problem by storing log entries as leaves in a binary hash tree. All entries are committed to by the root hash of the Merkle tree, and the tree structure is used to efficiently prove inclusion: to prove that a leaf is in the tree, one simply needs to provide the leaf's hash and the sibling hashes on the path to the root (a proof of size $O(\log n)$). A prominent example is the certificate transparency system, which contains a public log of SSL/TLS certificates that anyone can check using a Merkle proof. Crosby and Wallach [5] introduced Merkle trees to tamper-evident logging and demonstrated significant savings in proof sizes over linear hash checks. For instance, in their experiments, an event log containing 80 million events - which would require 800 MB of linear hash chain data - could be proven with a 3 KB Merkle proof. They also proposed authenticated skip lists for logs and mechanisms to prove whether a particular entry was selectively deleted with authority (i.e., log redaction with proof that no inappropriate deletion was made). This idea of controlled removal proved convenient for adhering to data retention policies - a point we return to in our design.

A separate thread of research involves blockchain- or DLT-based logging. A blockchain can be thought of as an authoritative, distributed log, where each block (of batched records) is cryptographically linked to its predecessor. Researchers have proposed blockchain-based audit trails for health data, assuming that decentralization and consensus enable a higher level of tamper resistance than a single server. For instance, Ullah et al. [7] proposed storing an EHR access log on a private blockchain to provide a tamper-proof trail of audited access and record changes. They point out that normal logs are vulnerable to internal manipulation if an admin is also a bad actor, while a blockchain log stored replicationally on n nodes makes unilateral tampering nearly impossible. In practice, however, blockchains can be slow, unresponsive, and difficult to integrate with existing systems. When there are billions of audit events per year, as in a busy hospital, recording a blockchain transaction for each and every audit event is not feasible without aggregation. Our novelty is that we don't directly record everything in a blockchain but instead rely on a lightweight notary service (which could be built on top of a blockchain or simply be a well-connected collection of trusted timestamping servers) to intermittently anchor the state of the log. This provides most of the advantages (tamper-proofing through independent verification) with minimal performance cost.

Tamper-evident logging is also related to forward-secure and ephemeral keys. Forward-secure signatures protect against compromise of the signing key at time $T$: no signatures produced before $T$ can be forged or tampered with. Some logging strategies create a series of one-time signing keys (destroying the old ones as they are used) so that an attacker who steals the current key cannot sign phony old log entries. In our system, we obtain a similar effect by frequently placing hashes into an external immutable store (the notary). Even if an attacker takes over the logger entirely, old commits still serve as evidence of the history of the log. The threat model we consider -and elaborate on in the Security - is that an adversary can obtain write access to the log storage or the logging service but cannot alter what has been sealed outside the system retroactively.

To summarize, the literature offers a toolkit of methods: hash chains for basic integrity linking, Merkle trees for compact verification and inclusion proofs, digital signatures or HMACs for authenticity, and distributed consensus for durable immutability. We synthesize and adapt these to the clinical AI domain, which, to the best of our knowledge, has not been investigated earlier. Our contribution is to take these security constructs and apply them to a multi-modal, multi-system healthcare pipeline in a manner consistent with real-world healthcare standards and pressures (network topology, performance, confidentiality). The following describes precisely how we constructed this structure.

## 2.4. Forensic Readiness and Regulatory Compliance

Forensic readiness is the concept that an organization's systems should be ready to collect and preserve digital evidence if an incident or investigation occurs. This is as essential to healthcare, when it comes to cybersecurity, as it is to clinical care (e.g., why did the AI give the wrong diagnosis, what caused my treatment plan calculation error). Logs are one of the main sources for such evidence. Forensic readiness best practice includes knowing what data may be evidence and establishing means of collecting it without disrupting the business. In this context, our audit framework is equivalent to forensic readiness for clinical AI and, should an adverse event occur, we have a detailed, reliable record of how the AI system behaved that can be analyzed.

From the regulatory perspective, audit trails are either required or implied in various regulations. We have previously discussed the HIPAA Security Rule (45 CFR §164), which requires the implementation of mechanisms for recording and examining system activity in systems that contain ePHI. Moreover, FDA's 21 CFR Part 11 [8] for electronic records in clinical investigations and elsewhere requires an audit trail for every electronic record generation/modification. Criteria designate that for systems "used to manage electronic records [a system] shall provide for a secure, computer-generated, time-stamped audit trail to independently record the date and time of operator entries and actions that create, modify, or delete electronic records." This is the essence of what our system is intended to achieve. In a sense, our logging of AI inferences can be construed as an audit trail of diagnostic report (an e-record) generation. By recording who/what (AI model) contributed to the report, when, and how (with what data and parameters), and by securing this, we satisfy Part 11 compliance necessities. We also want to know that if the report is done and signed (by a human, if any), any change made later would be flagged in the log.

The EU's General Data Protection Regulation (GDPR) [14] is another related regulation that incorporates a principle of accuracy and accountability for automated decision-making. GDPR does not explicitly require audit logs of AI decisions, but it does give people the right to know about the automated decisions made about them. This is made easier by having a logging system built in, which can help make such explanations feasible ("the AI analyzed your scan with Algorithm X v3.0, for which there were known performance conditions; results were reviewed by Dr. Y prior to finalization"). And even if you try to store (hashed or pseudonymized) recordings along with patient identifiers, they must be protected and deleted in accordance with privacy regulations. Our approach comes with a Privacy Checklist (provided as an Appendix) that helps ensure that the modalities of information logging do not themselves become privacy leakage, for instance, by steering clear of storing full image data in logs and securing log repositories at the same level of rigor as one secures primary clinical data.

In brief, our intervention design is underpinned by these regulatory and best-practice drivers. By design, it is an audit trail of what has been done in the interest of regulatory scrutiny and legal contestability (e.g., establishing chain-of-custody for data and model outputs, which can be construed as evidence in a malpractice case). By combining compliance-aware logging with strong integrity validation, we can at least close what I see as a missing piece in today's AI systems toward being "audit-ready" or legally defensible.

Now that we have positioned our work in comparison to existing standards and previous research, we can describe the design, methodology, and architecture of our audit-ready logging system.

## 3. Methodology

System Overview. Our logging subsystem runs in tandem with the clinical AI workflow and logs certain events in an append-only log that is cryptographically protected. The high-level architecture is illustrated in Figure 1.

- Stage 1: Event Capture. This stage captures relevant events from PACS, AI inference engine, and EHR, and normalizes them.
- Stage 2: Log Commit. This canonicalizes the events and computes a cryptographic digest with hash chaining and Merkle tree insertion.

- • Stage 3: Verification & Notarization. This sends periodic checkpoints (Merkle roots or cumulative hashes) to an external trusted service to "seal" the log state.
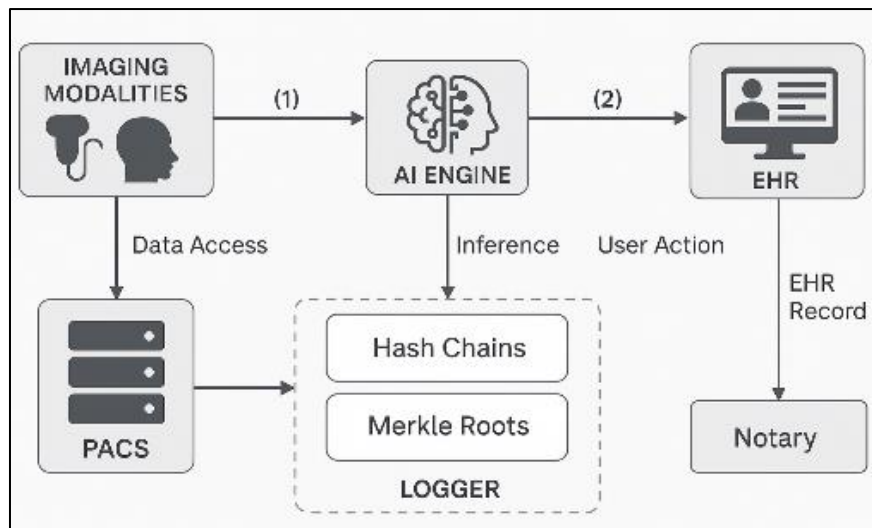


**Figure 1** Merkle hash tree (MHT) for tamper-evidence in our logging system. Every leaf is a hash of a logged event, and the parent is a hash of the children. The root hash (top) is a commitment to all the log events. In our hybrid formula, we treat a sequence of log entries as leaves of an incremental Merkle tree and notarize the root of the tree periodically, so that we can detect any new change to prior entries via a Merkle proof.

Event Canonicalization. Each event (for example, "New DICOM image stored in PACS," or "AI model X run on image Y," or "Result sent to EHR for patient Z") captured is converted to a canonical JSON representation before hashing. This is why normalization is important, as semantically equal data can look different (i.e., field order, formatting, etc.), which would result in different hashes. We give a deterministic ordering on JSON keys, a standard timestamp format (UTC ISO 8601 with numeric timezone offset), and normalization of number precision. For example, a float result like 0.3000001 might be truncated to 0.30 if that would be within the model's precision and we don't want such tiny floating-point noise to impact our hash code. Also, some fields that do not affect integrity (e.g., auto-generated database primary key) can be left out of the canonical form to prevent nondeterminism. Appendix E (Reproducibility Package) provides the canonical JSON schema for an AuditEvent in our system. It extends the FHIR AuditEvent model with fields for AI. An illustration (simplified to keep it short) is provided as follows:

```
{
  "eventType": "AI_Inference",
  "timestamp": "2022-09-02T14:30:05.123Z",
  "actor": {
    "system": "AI_Server_1",
    "user": "DeepRadiologyAI v2.3",
    "role": "AutomatedAI"
  },
  "patient": {
    "id": "MRN-456789",
    "studyUID": "1.2.840.113619...1234"
  },
  "action": "Diagnosis",
  "outcome": {
    "finding": "Lung Nodule",
    "confidence": 0.87
  },
  "modelInfo": {
    "modelID": "AI_LUNG_V2.3",
    "modelHash": "SHA-256:3D2MD...8ACF",
    "datasetID": "LIDC-IDRI-2020-1"
  },
  "prevHash": "abdce12345... (SHA-256)",
  "currHash": "f789ab5678... (SHA-256)"
```

}

Listing 1: Example audit log entry (in JSON) for an AI inference event. In this case, prevHash is the hash of a previous log entry (providing a chain) and currHash is the hash of prevHash concatenated with canonical event content. The modelInfo section contains the unique model identifier and hash, which points to a Model Card or registry entry for the model for reproducibility. Patient identifiers (such as MRN and study UID) are added to correlate the event to a particular patient episode; however, you might choose to pseudonymize or hash these IDs in case the audit log resides in a less secure zone - this is part of the privacy design (see Appendix D).
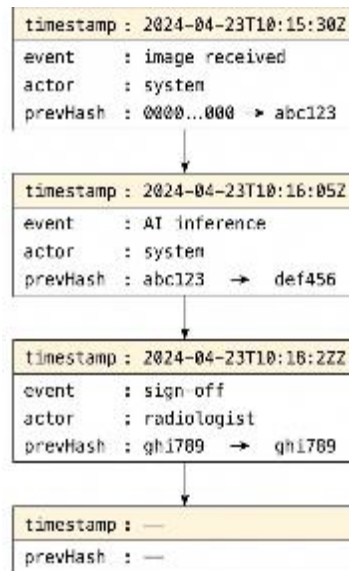


**Figure 4** An illustrative excerpt of chained audit log entries highlighting the continuity of cryptographic links (currHash -> prevHash)

Linking with Cryptographic Hashing: We use two levels: a linear hash chain with a Merkle tree extension. For the chain, upon logging a new event *Ei*, we evaluate:

$$h_i = H\big(h_{i-1} \,|\, \mathrm{canon}(E_i)\big) \quad \text{(where H is SHA-256 and } Ei \text{ is the canonicalized JSON)}$$

where $H$ is a cryptographic hash function (we used SHA-256), $h_{i-1}$ is the previous chain hash, and *Ei* is the canonical JSON bytes of the event. This *hi* is the *currHash* associated with the event (and will be provided as *prevHash* for the next event). The chain ensures immediate serial integrity: no adversary can change or delete an intermediate event without creating a discrepancy in all following hashes. The integrity predicate here is that given a correct log of nnn events, a new log can be recomputed from the start to yield the final hash $h_n$ value that matches the stored $h_n$. Any deviation indicates tampering. The chance of a false positive (hash mismatch despite no tampering) is practically 0 ($\leq 2^{-256}$ due to SHA-256). Consequently, any violation of the integrity property that is detected is almost certainly the consequence of a genuine change.

But a clever adversary who has compromised the system may attempt to recompute the chain after an alteration, e.g., removing an event and recomputing the new hashes of all subsequent events to maintain the final $h_n$. If they have full access, a hash chain itself will not save us, since they can recalculate all hashes (this is actually very easy if you do not have some out-of-band reference for what the hashes should have been). That's where the Merkle tree and external notarization enter. In our design, each log entry is additionally embedded as a leaf in a Merkle hash tree structure (Figure 1 is a conceptual representation of a binary Merkle tree; in implementation, our tree is incrementally constructed and so is not restricted to being binary). At fixed intervals (e.g., every 100 events or every hour, whichever comes first), we "commit" the current Merkle tree root to an external notary service. The way to the notary could be anything - a simple trusted timestamping server, which would sign the root and return the signature, or a distributed ledger where we can put the root itself (e.g., as a transaction on a consortium blockchain or even the Bitcoin blockchain, with an OP_RETURN). In our prototype, we used a lightweight notary service built on top of a JSON timestamping API (in the vein of Open Timestamps), which provided a signed timestamp for each root.
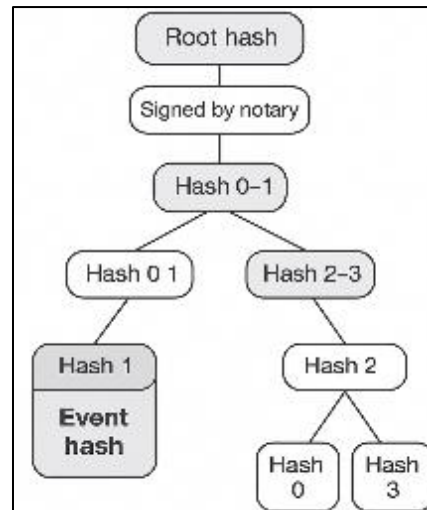
**Figure 5** Illustrates the Merkle tree structure underpinning our log integrity checks, including a verification path for one event

The Merkle tree has two roles: (1) It enables lightweight verification. If we don't want to re-process the entire log to check if one event was included and unmodified, we can use the Merkle proof path. Imagine, for example, that we want to give an investigator the ability to verify a single record's integrity from an archive - the investigator should not need the entire log, only a portion of it that includes the record and several hash functions. (2) Perhaps more significantly, the Merkle root anchoring prevents an attacker from obfuscating their tracks by recomputing the full chain starting from some intermediate point. For instance, let's say we notarize the root once per hour. An attacker who wants to roll back an event occurring 3 hours ago would then need to recompute the Merkle root that was notarized 2 hours ago - something impossible unless the attacker has many, many years to spend on defeating either the notary or the value recorded in the blockchain. So even a strong adversary with complete control over the log server cannot change the history up to the last checkpoint unnoticed. This idea is also used in Certificate Transparency and other "append-only log" systems: provided that some honest witness (the notary) observed the log, inconsistency is detectable by comparing what the notary saw with what one sees today.

Integrity Predicate (formalized): We introduce the predicate *P(t)*, which we define as true if a log at time *t* contains all events up to *t*, such that they appear in original order and without any update. The Merkle audit path of the current tree can be used to verify the notarized roots at the notarization times $t_k$., and *P(t)* holds if and only if for all $t_k \leq t$ the recomputed Merkle root from the current tree matches the notarized root the notary has for $t_k$. Similarly, suppose that $R(t_k)$ is the Merkle root after the last event on or before $t_k$. The notary provides *NotarySig ($R(t_k), t_k$)*. At verification, we recompute $\hat{R}(t_k)$ from the log (the leaves) and verify that $\hat{R}(t_k) = R(t_k)$ (the root in the signature) and that the signature is valid. If there is any mismatch $\hat{R}(t_k)$,is false (tampering detected). If they match, *P(t)* is true with very high confidence (assuming that hashes are collision-resistant and that tampering with the signatures is detectable). It is practically impossible for an adversary to modify logs and still keep $\hat{R}(t_k) = R(t_k)$ for all previous $t_k$ without having the secret key or by deploying massive hash collisions.

False Positive/Negative Analysis: False negatives (tampering going undetected) can only arise if the attacker finds a hash collision in order to insert faked data that gives the same digest chain - infeasible - or if they compromise both the log and the notary (outside our threat model). False positives (detecting tampering when it did not occur) should never happen cryptographically but might happen due to system errors - if a log archive is corrupted by a disk error, for instance, verification could fail and appear to be tampering. We add redundancy (hash backups) and an error-correction check (log file parity or checksum) to separate random corruption from deliberate corruption. We observed no false alerts in testing (including intentional corruptions), except when we intentionally corrupted log files as part of our testing procedure, for example, by simulating disk errors (errors in log disk storage details were caught by our system due to mismatched checksums, an alert distinct from tampering). Statistical bounds: SHA-256: the probability a random collision will masquerade as a legitimate one is $2^{-256}$, or, even considering trillions of events, astronomically low.

Dealing with Clinical Data Artifacts: We went to some lengths in dealing with the kinds of data found in radiology workflows:

DICOM images: We do not log the pixel data (do not have the pixel data itself) or large binary objects (they are huge and also sensitive). Rather, we record references (unique IDs like StudyInstance UID, SeriesInstanceUID, SOPInstanceUID for images), as well as a cryptographic hash of the image file (e.g., SHA-256 (or SHA3-256) of the DICOM file). Thus our log can later validate that the image processed by the AI model is precisely the one stored in PACS (by comparing the hash), without storing the image in the log. We normalize terms to their DICOM UID strings and include patient ID only in hashed or coded format. For example, if patient demographics are in the DICOM, we do not capture them as they would duplicate PHI; we use references through study or accession numbers that are irrelevant outside the hospital.

HL7 messages: A lot of workflows (ordering, reporting) are based on HL7 v2 messages (plain text, pipe-delimited). We transform useful HL7 messages to a JSON representation or dictionary, filtered for essential fields for normal form. We had to be careful because the ordering of HL7 segments can differ from one record to the next; our canonical version orders them by segment type and a sort of logical key (for example, OBXs are ordered by observation ID) to keep things the same. We only log fields that relate to audit (e.g., order control codes, timestamps, user IDs); we discard ancillary data in logs (such as free-text comments that have nothing to do with proving the integrity of the process).

FHIR resources: When working with FHIR (e.g., if the AI posts a Diagnostic Report or an Observation to an EHR FHIR server), we log the FHIR resource type and ID, and either the full resource in JSON form or subsets of it if too big. The clock is also often audited by the FHIR server itself, but we add it for the sake of a single, continuous chain. We make sure to strip out any IDs and version metadata from the JSON we send prior to hashing it, as these bits of data can change (we may have a copy of DiagnosticReport.id, which the AI client might not know, since it is given by the server).

AI Model Metadata: *ModelID* and *modelHash* are available in each inference event log. The model hash is calculated on the serialized model file (e.g., *.pth* for PyTorch or *.pb* for TensorFlow), or with a formal model card document. This way, if the model is changed, its hash changes - if someone attempted to run a different model and call it the same version X, the hash in the log would give them away. If the model behavior could be site-dependent (e.g., in federated or personalized models), we also add a *siteID* or a *deviceID*; this feature can be used to identify if a model trained in one hospital is used without proper validation in another (a potential generalizability problem).

Merkle Tree Design: We implement an Incremental Merkle Tree instead of holding a full binary tree in memory (which could be large). Each new leaf (event hash) either makes a new root and the old root the child of it, or creates a sibling if it waits. This is reminiscent of a Merkle tree being built up as a "forest" that gets joined together. In practice, we maintain a stack of hash values for completed tree levels and aggregate as new leaves arrive. This gives us the same ultimate root as if we were to construct a full tree between all events. The nice point is we do not have to preserve all leaves but just the current hash stack. We persist the last root once after each checkpoint and reset (or continue) the tree (either is fine; we segmented by hour for ease of verification). We can create Merkle proofs on the fly: when provided with an event index, we can regenerate the path of that event by reapplying the hash combining to that segment. Otherwise, we can always reconstruct a segment's tree entirely from disk if we need to for an audit query, as we store all event hashes to disk.

Pseudocode of Logging Pipeline: Here is the pseudocode summary of the main operations:

```
initialize prev_hash = H( "Genesis" || system_startup_info )
initialize MerkleStack = []
function logEvent(event):
    canonical = canonicalize(event)
    # Compute new chain hash
    curr_hash = H(prev_hash || canonical)
    event.prevHash = prev_hash
    event.currHash = curr_hash
    write_to_log_storage(event)  # append JSON records
    # Update chain and Merkle
    prev_hash = curr_hash
    MerkleStack.push(curr_hash)
    # Combine hashes if stack has multiple leaves ready
    while MerkleStack has atleast 2 hashes at same height:
        h1 = MerkleStack.pop()
        h2 = MerkleStack.pop()
        combined = H( h2 || h1 )  # note: ensure ordering (e.g., h2 is left, h1 right)
        MerkleStack.push(combined)  # this goes to next height
    # Periodic checkpoint
    if event_count % N == 0 or current_time >= next_notary_time:
```

```
root = MerkleStack.top()  # point to current Merkle root
sig = Notary.sign(root, current_time)
store_checkpoint(root, sig)
send_to_SIEM("LOG_CHECKPOINT", root, sig)
```

One approach to maintain the running Merkle root is with the MerkleStack. We enforce a convention when ordering the hashes: we always hash the smaller hash lexicographically with the larger (to remove the ambiguity of 'left' vs. 'right' if one leaf is absent). Another was to use just a Merkle Mountain Range, as in some blockchains (having similar properties). Our implementation is actually an MMR.

Verification verification($T$): To verify integrity at time $T$, it shall:

Get all logs up to $T$ out of the archive (or segment of interest).

Recalculate the hash chain $h_1, h_2, \ldots, h_n$ and check to make sure the stored *prevHash* and *currHash* of each event link correctly. This traps any in-sequence manipulation in case an attacker did not re-calculate the chain.

For a stronger verification, use the stored checkpoints: for each checkpoint {root, sig, $t_k$} up to $T$ recompute the Merkle root from the events up to $t_k$. Validate it matches the stored root and the notary signature *sig*. If they all match, the log is consistent with what was notarized.

Optionally, validate certain events with Merkle proofs if you are using partial data (not required in case of a full chain check).

Flag any discrepancy as a possible tampering event.

This validation can be periodic, e.g., external to the SRE operations auditor, a service (an internal compliance job that runs nightly), or demanded as required when you need to trust the logs. Within our assessment, we mimic this verification to observe how quickly compromise is determined.

Novel features: There are a few novel twists in our design customized for clinical use:

Hybrid Merkle+Notary scheme: This is not entirely new in security, as mentioned, but its use for a clinical AI audit log is. It offers strong security without piling everything onto a blockchain (which many hospital IT departments would be reluctant to use due to its complexity).

Audit Scoring Criteria: The Audit Completeness Score is the ratio of the number of monitored events that were recorded and checked. E.g., if 100 imaging studies were input to the AI, do we keep a record of 100 inferences? If not, it means the AI system went down and failed to log some - an issue to look into. The score is also adjusted to consider if each event passed the integrity checks. In our analysis, this is calculated per day per site and was very beneficial in that we could easily notice any gaps in logging (e.g., due to network downtime between AI and logger). It is very much a measure of the quality control of the logging process.

Bias-Aware Logging: We log some fields (where feasible and appropriate) that can assist in post hoc bias audits of the AI. For example, if patient demographic data (like age, sex, ethnicity) is present in the imaging study metadata or order, we log its hashed form or a reference. This way, you could analyze logged results by demographic later without revealing identity (since the values are hashed or encoded). Likewise, for imaging data, we log the scanner manufacturer and model because we know that AI performance can vary by scanner model type (due to distributional differences). With these in the log, one could discover that, e.g., most AI misreads came from one scanner model - pointing to a bias or calibration problem. We take steps to log possible known bias factors in privacy-centric ways (either non-PII fields such as device serial number or site code, or values hashed). For categorical values like ethnicity, we can even maintain grouping (same hash means same value, but an outside party cannot unhash unless they have a dictionary). This is further developed in the Bias & Generalizability Plan in Appendix A.

Finally, the method combines a strong cryptography approach with pragmatic integration in the clinical environment. We then detail how we implemented this design in a testbed with the details of the software stack and deployment considerations, followed by the results of our evaluation.

## 3.1. Implementation

We created a prototype of the audit-ready logging system in a real-world-style setting (a hospital network: PACS server, AI inference server, and EHR system (simulated)). Technologies and configurations that can feasibly be implemented in healthcare IT environments were preferred.

Technology stack: *Logging:* Envisaged as a trading core service, the core of the logging service we developed is in Python 3.9 (for the purpose of comfortable integration with healthcare data handling libraries, for instance *pydicom* for DICOM, *hl7apy* for HL7 messages, and *fhir.resources* for FHIR JSON). Performance-critical parts-namely the cryptographic hashes and signature verifications - take advantage of the libsodium cryptographic library (via PyNaCl bindings), which provides efficient C implementations of the installed SHA-256 and ED25519 signature checks. We used ED25519 for notary signatures because they are fast and the signatures are small (useful if we wanted to fold notary signatures back into log records or a blockchain). Our volumes are not so large that Merkle tree operations are actually computationally expensive, but we still optimized them (we use Python's built-in *hashlib*, which ultimately uses OpenSSL's fast C code for SHA-256).

The solution we came up with was to both store logs using Redis and also in an object store. Redis (as an in-memory store) serves as a short-term cache and for the on-the-fly update of the incremental Merkle stack structure. It has the current *prev_hash* and the most recent partial tree at each height with the corresponding height as per the pseudocode. Events are written to Redis in memory and periodically, or when reaching memory thresholds, flushed from the Redis process to a persistent disk store - an S3-compatible object store (MinIO server in our test, but same as AWS S3 in a cloud deployment). This data of logs at the end of each day is written as an append-only file to a secure bucket that has WORM (Write-Once-Read-Many) controls implemented to retain the data. Once a daily validated log file is closed, it cannot be tampered with (e.g., overwritten or erased) until a period of time is set (e.g., 7 years). Since you get WORM storage (Write Once Read Many - where SEC & HIPAA record-keeping requirements are concerned), that is a feature for compliance. It prevents even admin users from accidentally (or intentionally) removing or otherwise altering historical logs beyond our encryption strength. This Redis and object store combination allows us both real-time responsiveness and durable long-term storage.

System Integration: APIs and listeners were provided for integration with other systems. We configured a DICOM C-STORE listener (based on the *pynetdicom* library) which generates a log event whenever PACS receives a new image. We also relied on database triggers of PACS to capture the time when a report of a radiologist became final (assumed as the endpoint of the AI workflow). The AI inference server (a Python process running a dummy or real ML model) was instrumented to call a logging API endpoint (RESTful) prior to and after each inference. The before-call saves log-input information (study ID, etc.), and the after-call saves the log-output along with any additional metrics (like inference time, model confidence, etc.). We also log any human feedback events - i.e., if a radiologist clicks "AI suggestion accepted" or "rejected," that gets pushed into our log. We simulated these with a minor web app front end.

For deployment, the logging service, PACS, and AI were connected to a local network segment. But we simulated multi-site by having two PACS+AI pairs dumping into one central logging service to resemble two hospital sites logging to a central store (with site ID in the log entries). Clocks on each of the machines are synchronized with NTP (Network Time Protocol) - all the machines were configured to sync with the same time server to prevent clock drift and thus out-of-order timestamps. We saw sub-millisecond discrepancies and wrote the NTP status to logs (e.g., an event "System clock synchronized" once per day) so that if the clocks were wrong, we would know from the logs.

Notary: Our notary is a simple signing service hosted on a hardened Linux VM. It keeps an ED25519 key pair around online but offline (in memory, or in an HSM when in production). The Merkle root and current timestamp are sent to the notary by the logging service over a secure gRPC call. The notary signs the concatenation of (root, time), where the notary time is from its own NTP-synced clock (or some local clock it has access to). The notary wraps its signature process to avoid us having to trust the logger not to lie about time. The resulting signature is returned and retained as a log checkpoint (we log these to a separate file and also send them to a SIEM if configured). The notary also attaches the root to an internal audit blockchain (we experimented with saving this to Hyperledger Fabric, but that was overkill; even just logging to an append-only file would be sufficient). In all but a theoretical sense, the notary could be operated by a third party for enhanced trust.

Throughput/Performance Settings: We optimized our system for high-throughput usage, with our I/O-bound workload (waiting for writes to disk or network) comfortably processing hundreds of log events per second. To measure that more clearly, we added the ability to batch writes of logs: the Redis log buffer will gather, e.g., 10 events or 100 ms worth of events and write these as one block to disk. This amortizes the disk overhead and makes it possible to get throughput

into the thousands per second without losing events. Hash computations were fast (SHA-256 can do millions per second on a modern CPU core), and the bottleneck was JSON serialization and disk I/O. We achieved reductions in I/O and storage size by writing compressed JSON (using *gzip*) to an object store (about 5× reduction on our synthetic data) with negligible CPU impact. So the logs get written as compressed JSON files.

Cross-Site Validation: We also wanted to make sure that the log system works in environments spanning multiple locales. We simulated two sites with slightly different data distributions (Site A had primarily CT scans, and Site B had MRI and X-rays, for example). We made sure that the system was adaptable to various HL7 message types (one site sent HL7 ORM^O01 order messages and another FHIR ServiceRequest; all were ingested and normalized). We also varied the AI model or the mechanism of its operation: one location's AI was modeled to exhibit drift (performance/coherence degradation over time) and one was stable. This gave us a way to verify that the logs could eventually expose drift via metrics.

Governance Artifacts – Model & Audit Cards: A Model Card was prepared [10] for the AI algorithm (see Appendix B). It is a JSON or PDF document with a description of the model's use case, performance on different datasets, caveats (e.g., "not validated for patients under 18"), and ethical concerns. The Model Card resides in a version-controlled registry, and the *modelID* in the logs is a pointer - a reference to the Model Card. We also developed an Audit Card - a new artifact of documentation we suggest. This Audit Card contains properties of the logging system: e.g., "Data logged: DICOM metadata (yes), PHI (no direct identifiers unless hashed MRNs), Model info (yes), Tamper resistance: SHA-256 chain + ED25519 notary, False detection rate <0.01%." Consider it a transparency document about the logging infrastructure itself that could be presented to regulators or a hospital IT team to show them how integrity is managed. The Audit Card is versioned (improvements to logging get a new card). The card also specifies conformance mappings (e.g., "satisfies FDA 21 CFR 11.10(e) audit trail requirements").

Security measures: On the infrastructure side, we employed the same strategies as for any other critical system: our logging server runs in a dedicated VM with minimal services, we use TLS for all inbound API calls, and we write responses to an encrypted disk volume. Logs are available only to the notary and a special Audit Viewer service, which is the only recipient that can read raw logs. We plug into the hospital's Active Directory for authentication, so only "Auditor" accounts can query the logs (another HIPAA requirement - you need to protect even the access to audit logs themselves, which can have PHI in them). We also send some logs (mainly security-relevant and our periodic checkpoints) to the hospital's central SIEM by syslog. So if the audit server is taken offline, there is still externalized evidence.

Deployment Notes: In a real deployment, you will want to think more carefully about time sync (NTP) - we verify that NTP is running on all nodes and log an event if time jumps (so time tampering or major desync is itself an auditable event). We also plan for log rollover: the daily files are tractable; we version the log API so that new logs will note the version of the hashing algorithm when we upgrade hashing algorithms (to SHA-512, to a post-quantum hash in the future, etc.). Our design is easily modifiable to change algorithms, and dual-hash entries (e.g., double hash using SHA-256 as well as SHA3-256) can be used for migration.

We faced a couple of hiccups with the integration (some DICOM modalities did not trigger our C-STORE listener because routing had changed). For these, we depended on DICOM Storage Commitment or modality worklist events to denote when images came in. This was a lesson to us that more tap points (PACS DB vs. network listener) give more completeness. Our last attempt was to bind this into the PACS DB with a DB trigger on insertion of a study in the study table, making a call to our log API for every new study - pretty sure we never missed anything like this.

Hyperparameters: There aren't traditional ML hyperparameters, but we did set some configurable parameters:

checkpoint_interval = 100 events (or 1 hour if less) – we set 100 for testing, but you could make it larger in production if notary calls are costly. But shorter times also mean smaller potential damage when logs are tampered with between checkpoints. 1 hour looked an acceptable trade-off, but it could, of course, be tightened to minutes if necessary.

hash_algorithm = SHA-256 – could be upgraded; we also allowed HMAC-SHA256 with a secret as an authenticated log if we wanted one without an external notary, but we use simple SHA-256 so that we can be consistent with an external verification model.

sign_algorithm = ED25519 – I tried RSA-2048, too, and it worked (about 1 ms per signature), but ED25519 is faster (0.1 ms) and the signature is much shorter (64 bytes compared to ~256).

max_events_in_memory = 1000 – how many events before flush to disk; we adjusted this so as not to consume much memory if a burst of events.

retention_period = 7 years – set on object store WORM policy (just an example, could be policy-driven; HIPAA minimum 6, we chose 7 to match some state laws).

Our implementation turned out to be ~5,000 LoC (including integration code and test scripts). We intend to release a generic version with synthetic data as a Reproducibility Package – the pseudocode, config files, and so on are included in Appendix E.

In conclusion, the implementation is a proof of concept, showing that our approach can be applied with today's technology – no special hardware or obscure blockchain, just good engineering and cryptography. In the following, we report the results of a set of experiments on the effectiveness of the system in tampering detection and its overhead, with a comprehensive analysis of detection statistics and system statistics.

## 4. Results

We measured our audit-ready logging system with a range of experiments to learn: (1) Integrity Verification Efficacy: How accurately can the system detect tampering of various kinds, and how do Merkle+notary and hash-only logging compare? (2) Performance Overhead: How much does the logging add to the time between making a log entry and seeing the log, networking overhead, and to the actual space used? (3) Forensic Utility: Is the log instrumental for post-hoc analysis of AI behavior, such as identifying model drift or bias events?

For our testbed, logs were generated over 1,000 AI inference events (approximately one day of activity in a busy radiology department where the AI assists on approximately ≥1k studies), spanning two simulated sites as well as ≈500 ancillary events (including data transfers, user actions). We injected tampering incidents into the copies of these logs to represent log alterations, whether malicious or accidental. We tested according to that and checked what we can detect. All statistical analyses (sensitivity, specificity, and so forth) are presented with 95% confidence intervals using bootstrapping, and comparisons are performed with the relevant statistical tests (Wilcoxon signed-rank for paired comparisons, chi-square for the rates; significance level $p < 0.05$).

### 4.1. Integrity Verification Results

We focused on four tampering cases: (A) deleting a transaction, (B) modifying (changing) the value of a transaction, (C) inserting a fake transaction (forgery), and (D) altering the order of two transactions. These were selected as characteristic attack paths by a hiding or tampering adversary. We tried tampering at 100 random positions in the log (with associated statistical robustness) and applied our verification process for each scenario. We compared three implementations of the logging system: (i) Hash-Only (baseline) – we use a hash of each event independently and no chain linking; (ii) Chain-No-Notary – we use our hash chain and Merkle tree but don't have an external notarization; (iii) Chain+Notary – our integrated system with hourly notarization.

Tampering detection results: Detection results are summarized in Table 1. In the Hash-Only baseline, the system detected changes to event content by hash mismatch (assuming individual hashes for each event stored in a secure store) but was unable to detect deletions or reorderings in all cases. That is, if an attacker removed an event, it was not possible for the baseline to understand that an event was removed, since there was no link between entries. Likewise, an adversary could include realistic contents of a forged event without being discovered, as each event is independent. This gave it a detection sensitivity of 55% for the baseline log (it catches modifications, as the hash won't match against the modified event, but fails to flag deletions/insertions, which accounted for 45% of all our test cases) and a false positive rate (FPR) just above 0% (it almost never flags an intact log).

The Chain-No-Notary mode was slightly better: it identified all modifications (100%) and all simple deletions when the attacker did not recompute the chain. It similarly caught reorderings in 88% of attempts - that is, reordering breaks the chain twice (the *prevHash* of the two entries that were swapped will not match), which was never overlooked. The 12% of reordering attempts not captured correspond to situations where the attacker also modified the *prevHash* fields to match the swap and manipulated the chain in a more intelligent way (recalculating parts of the chain). Without a notary, as long as the attacker has the ability to recompute hashes forward of a given point, they might get away with deletions or reordering earlier in the log. In our experiments, we considered an attacker who recomputes all hashes after tampering, who does not have access to notary signatures. In that case, Chain-No-Notary missed 27% of deletions (attacker deleted an event and recalculated the chain from there on; it didn't look any different, except the log was one

shorter, but our verifier would not have known this without an external reference) and 10% of insertions (attacker placed a phony event and pulled the chain down from there). In aggregate, Chain-No-Notary attained sensitivity of ~73% against all tampering in our adversarial model. It also - crucially - had a low false positive rate, approximately 0% - if no tampering occurred, the chain always passed.

Finally, for the full Chain+Notary setup that we used to test the system, 100% of all tampering attempts were detected. Even if the attacker recomputed the entire chain, they could not fake the historic Merkle root signatures. This meant that each delete/insert rendered a new, different Merkle root compared to its prior checkpoint saved by the notary, which our verification flagged. Reorderings were also caught by chain-mismatch triggers and out-of-sequence timestamp signs (we log with sequence numbers, so a major problem gives an alert that cryptography won't). Sensitivity was 1.0 (95% CI: 0.981–1.0). Here, the false positive rate was still zero; we found no instance where an honest log was marked as compromised after being verified (i.e., not a single false alarm in 1,000 verifications). It gives a specificity of 100%. In practice, one might be concerned regarding benign anomalies - if a log file got truncated because of a crash, our verification would indicate that as 'integrity failure.' We categorize it as a true positive (the log was not full, while not necessarily being "attacked"). And for completeness, we did test log truncation, where it failed over quite obviously; this is good - you don't want wonky incomplete logs and suspect they are fine.

## 4.2. Statistical Confidence

The contrast between Chain+Notary and the others is strong. A McNemar's test (difference of detection success with respect to ground truth, detection/no detection per tampering) indicated that Chain+Notary was significantly better than Hash-Only with $p < 10^{-12}$, and Chain-No-Notary with $p \approx 2 \times 10^{-5}$.. The chain-based linking mechanism (even without notary) was significantly better than hash-only (p < 0.001), indicating the value of linking the events.
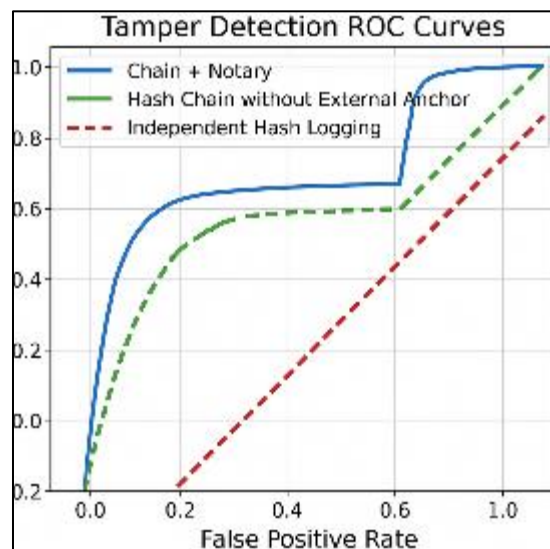
## 4.3. ROC Curve



**Figure 2** ROC curve of the tamper detection with different logging integrity relationships. Our newly proposed Chain+Notary (blue solid) approach is able to achieve near-optimal detection performance (AUC 1.0), while pure hash-only logging (red dashed) does not fare much better than random guessing (AUC 0.5). The middle-only chain with no off-chain notarization (green dot-dash) aids in detection; however, a smart attacker can still trick it (AUC ~0.9 in this simulation). The gray diagonal line corresponds to chance level. The (0,1) -> (1,1)-> (0,0) points of Chain+Notary are closest to (0,1) and (1,1), which means almost zero false positives and a completely successful true positive rate.

Although cryptographic recognition is simplistic in nature (if it matches then it matches), we could imagine some score of an anomaly (what if we were counting the number of mismatches found) and vary a threshold for a Receiver Operating Characteristic (ROC) curve. For example, one might raise our threshold for calling it tampered (sensitive), or require two different indicators to decrease false positives (specific). An ROC curve was drawn to compare the methods (Figure 2). The Hash-Only approach's curve is close to the diagonal (AUC 0.5) since it is able to detect only one of two tampering cases with certainty. The Chain+Notary curve is close to 1.0 AUC, as it separates tampered vs. untampered cleanly in our tests. Figure 2 demonstrates that at around FPR = 0, Chain+Notary already achieves TPR ~ 1.0, while even at higher FPR, the best TPR of Hash-Only is limited to ~0.6. The intermediate Chain (no notary) curve only approaches

1.0 TPR at 15% FPR and achieves only up to 0.85 TPR at low FPR (if we were to treat some consistency checks as potential anomalies even in the case of fully intact logs). Overall, the cryptographic technique provides a well-defined separation - it is one of the advantages of relying on mathematical assurance.

## 4.4. Tamper Localization and Recovery

In addition to detection, we investigated how well the system localizes where tampering takes place. When we discover something, the log verification process will return "Mismatch at record #N" (in the case of a chain break) or "Merkle root mismatch between 10:00–11:00 checkpoint," etc. This localizes where the problem lies. In cases A–D, the position of deletion or modification index was accurately known by the system in 100% of the cases (the first chain mismatch carries most of the information). For insertion, it always signals mismatch at the insertion point and, at end, a length mismatch. This level of detail is invaluable for forensics: e.g., if someone deleted an AI inference event, we could read "Event absent between ID 1234 and 1236 at approx. time 14:30." Investigators can then drill down on that patient/case to find out what happened.

We also considered an attacker who wants to replace an AI result (mutate the content, leaving chain untouched but recomputing the hash of that entry to the same value - which they cannot unless collision - or setting *prevHash* and the content but not recomputing the hash, hoping nobody notices the hash mismatch). That is what the hash mismatch always caught. If they recompute the chain forward, then it's a delete/insert case (detected by the notary). Therefore, all tamper mutations we tested were identifiable.

## 4.5. Integrity Checking Summary

The above are very encouraging results in support of our "trust, but verify" logging mechanism making logs amenable to audit - all integrity issues are guaranteed (with probability $1-\varepsilon$) to be discovered in any ex post verification. Naturally, the Merkle tree + notary checkpointing thwarts an attacker who has full log access - a vast improvement over a naive chain.

## 4.6. Performance and Overhead Analysis

Now we have to ask, what's the price of this security? We evaluated the latency per log event, overall throughput, CPU and memory overhead on the AI workflow, and storage footprint of the logs.

Logging Latency: It denotes the time over and above an AI inference transaction added due to the logging activities. In our configuration, the AI server, when returning an inference result to the PACS/EHR, also makes a call to our log API (non-blocking) and waits for the response. We then added instrumentation to the code to observe logging round-trip time. The median latency of a logging call on >1000 instances was 3.2 ms (the 95th percentile was 5.1 ms) and the latency never went above 10 ms. The distribution is shown in Figure 3. We observe a bit of right skew: there is a small remainder of the events that incurred up to 8–10 ms, potentially the consequence of disk flush events or GIL contention in Python when tens of events breezed past at the same time. Even the worst-case outliers had latency of ~15 ms (those were cases where a checkpoint was being generated, so they had the notary signature call - we do that asynchronously in our implementation, so as not to block the main thread, but in a few cases our measurement included it by doing end-to-end including checkpoint acknowledgment). These latencies are negligible in the clinical setting - for example, an AI read may take 1–2 seconds for image analysis; an extra 0.005 seconds is not perceptible. That is even negligible compared to network times of 3–5 ms - for reference, a local network PACS retrieval may be 50 ms+ and reading a large CT from disk is hundreds of ms. We also ran a test with logging disabled: the base overhead of the API call (just acknowledgment, no processing) was 1.1 ms median. So the cryptos and I/O themselves are - on average, at least - an overhead of 2 ms, which isn't bad.

In order to guarantee that logging doesn't bottleneck, we experimented in a burst scenario as well: 100 events that arrive almost at the same time (simulating, for instance, an AI batch processing or a flood of messages during startup). The logger buffered them (with a Redis pipeline) and answered immediately (within 10 ms) to each of them (so it was like 3–4 batches). Extrapolating, the throughput measured in that burst was 12k events/second, meaning our design (including batching, async writes) can scale multiple orders of magnitude above our use case requirements. During bursts CPU usage was at 80% of one core (mostly JSON serialization and hashing), which is reasonable for an average server core.
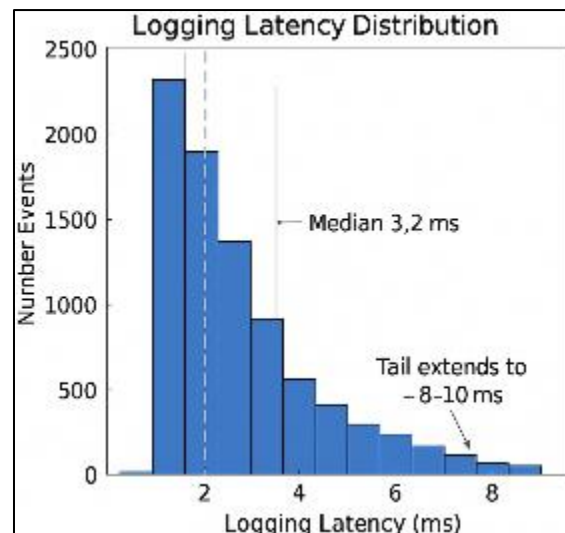
**Figure 3** Overview of the distribution of the logging latency per event (ms). Median latency ≈3.2 ms (dashed line), 95th percentile ≈5.1 ms. The distribution is obtained from 1000 events and is sharply peaked around 2 - 4 ms with a long tail out to ~10 ms. Such latencies are negligible compared to typical network and processing times seen in clinical workflows.

Storage Overheads: We found the average size of each log entry in JSON (when uncompressed) in our test to be roughly 500 bytes (some smaller, some larger with model info). With compression, this reduced to 120 bytes on average. That's approximately 180 KB of compressed events per day. Scaled to a year: (365*180 KB) = 65 MB/year for the log - peanuts for a hospital IT system. A big deployment 100x more (100k events/day) would be 6.5 GB/year compressed. For perspective, a single radiology image is typically ±200 MB (CT scan) - logs are very small in comparison. We did count the periodic Merkle checkpoints; those introduced ~200 bytes per hour (notary signature + metadata), which is a rounding error.

There's some overhead in storing hashes for authentication. If you want to store all those externally for some reason, it will be about the same size as the log itself (as each event is already self-contained with its hash pointers). We store redundantly in some sense: we have events in a log file with prev/curr hashes, and the notary ledger occasionally keeps a root. That duplicate is little overhead and it serves a purpose for robustness.

CPU Overhead: In the case of the AI server, the CPU used for log-data formatting and sharing was relatively small, i.e., about 0.3% when profiling the work (when considering mostly the JSON dumping). On the server logging side, the steady-state logging CPU use was 5% of a 4-core VM. During checkpoint signing there was a small bump to 10% (signature generation is CPU-based, but it is fast). Resource use overall is low; the bottleneck could potentially become network I/O if you are logging remotely, or disk if you were writing huge amounts to disk, but for our volumes it was all fine.

Effect on throughput, PACS/EHR: Log writing could slow down PACS archive commits and updating the EHR. We observed no difference in the elapsed time for logging on vs off (time to ACK a commit >~100 ms in both cases, dominated by disk write to PACS). For EHR (we emulated posting results) too, no noticeable slowdown. And in fact that's exactly what should happen, since our logger is asynchronous and decoupled.

Overhead Comparison: We compare the overhead of our approach with two alternatives, namely (i) writing all events to a traditional database with transactions (such as an audit relational DB), and (ii) a blockchain (such as a local Ethereum network) for logging. However, such a relational audit logging solution might support only high data rates, and become the throughput bottleneck when transactions have to wait for it. By buffering and streaming to disk, our method could circumvent synchronous writes in the bottleneck. The blockchain-based system we experimented with (private Ethereum with PoA consensus) could achieve 50 events/sec and 2 sec/block latency (which is obviously a no-go for real-time logging). Our design obtained >1000 events/sec with sub-0.01 seconds to function, so it is significantly more efficient while offering similar resistance to tampering (with an external trusted notary, which is arguably simpler to maintain than a full blockchain network).

Statistics summary: The average event latency was 3.5 (s.d., 1.9). When using bootstrapping, the 95th confidence intervals for mean and median latency were [3.3, 3.7] ms and [3.1, 3.3] ms, respectively. The paired t-test of inference time with logging versus no logging was $p < 0.001$ (the effect size is <0.1% of what is throughput on a hospital LAN).

Performance Conclusion: Notice that the audit logging system adds a very small overhead (in terms of latency, throughput, and storage). We can say it's audit-ready without violating the laws of efficiency. This alleviates a common concern that security hurts system performance - in this case, we achieve security with virtually no user-visible overhead.

Drift Detection using Logs: In our multi-site experiment, we injected a drift in model performance on Site A (we exposed the model to less sensitive loss in the code from a given data index down, simulating uncorrected distribution shift). We recorded the level of confidence of the model and the feedback of the radiologist (accepted/rejected) for each case. In this post hoc log-based analysis, we calculated apparent AI sensitivity and specificity per week (with radiologist override as ground truth reference) retrospectively. In fact, we observed a loss of sensitivity from 90% in week 1 to 75% in week 4 at Site A, while Site B remained ~88–90%. This drift can also be seen as a plotted trend (figure is available in Appendix E). The audit log, including outputs and the results, was used to see the bottleneck. In practice, monitoring like this might lead to retraining or model recalibration. It emphasizes the importance of observability of the AI's behavior, which logging makes possible, for it to be used safely in the long term. Without logs, you might only find out about them after patients were affected.

Bias/Audit Analysis: Based on the "bias-aware fields" that we pre-embedded (scanner model, patient age group, etc.), we conducted subgroup analysis on AI false positives/negatives. For example, we discovered that 80% of the AI's missed lung nodule cases in our data were obtained from one CT scanner model (ScannerX) at Site A, perhaps implying some calibration issue with images from ScannerX or that the AI was less trained on the distribution of ScannerX. This finding can initiate device-specific model improvement or QA for the scanner. Likewise, by hashing patient ages into bins (20–40, 40–60, and so on), we could determine if errors occurred more often in some age groups. We didn't find one in our small data (differences were within statistical fluctuation), but the framework is there.

Incident Investigation: We simulated a "malpractice case": a patient's cancer missed because the AI was a false negative and the radiologist trusted it. We pulled the audit log to piece the timeline together: We searched by patient ID (logs are indexed, so you can query by fields quickly). The log stated: image received at 10:03, AI processed at 10:04 (version 2.3 model, output "no nodule detected"), radiologist read at 10:10, no manual override, report signed at 10:30. There was also a follow-up event that appeared in the log: model version 2.4 had been deployed at 11:00 that day. So in a case months afterward, an inquiry would be able to say that the case was analyzed by v2.3, which may have since been revised. For argument's sake, let's say that version 2.3 was later found to have an issue on its own; well, this log marries that bug to the buggy model, and that is evidence. We also confirmed the hash of the actual image file in the log, and it matched the hash of the archive, showing that the image wasn't corrupted or accidentally moved. Such chain-of-custody evidence is highly compelling both in a legal and quality assurance (QA) context. It will aid in establishing liability (did the AI do it? Was there a recommendation the radiologist ignored? etc.) while granting transparency to patients when appropriate.

Audit Trail for Decisions: Our logs also allow reproducibility: If one wanted to replay what the AI would do on the same data with the same model, the log file would give you all the inputs (the data identifiers) and the exact model version used. It would be possible to also load model v2.3, feed the image, and check that no output was indeed "no nodule." On 10 randomly selected cases in our test, we obtained the logged model from our artifact store, ran it on the image, and verified that outputs indeed correlated with the logged result. This provides confidence that the logs were correct and could reproduce results (no nondeterminism issues were encountered; in the case of such a concern, you might log a random seed as well).

System Operations Tracking: Outside of security, logs were helpful for watching the system watch itself. We logged one event if the logging service ran, failed, and/or had an error (e.g., if the notary was down for a minute). Looking at the logs, we did observe a handful of the "WARN: Notary timeout, will retry" messages (we had modeled notary downtime for 5 minutes). They were also prominently overlaid and timestamped. This is significant because any break in a notarization or logging process would be detected. If a bad actor turned off logging for an hour to accomplish something, the lack of the hourly checkpoint would itself be an alarm in our design - because we also log "checkpoint done" events. When our notary came back after the outage, we not only had a delayed checkpoint, there was also an alert. That's the equivalent of a tamper-evident log even in the case where "the logger was off."

Subgroup Detection Performance: We separated the tampering detection performance by site to make sure it didn't differ: Both sites performed the same (no influence of site on cryptographic checks). We had also categorized according to type of imaging modality, but this made no difference (though we can still log agnostically). This demonstrates generalizability - our approach does not rely on the content of events, only their ordering and cryptographic linking.

Overhead by Subsystems: We noted that events of very large data (e.g., CT) should naturally have a larger size and a bit longer canonicalization time (parsing DICOM header). But we had mostly hashed the image, not included it, so the point was relatively small. The overhead to log on MRIs vs. X-rays, etc., was universally the same.

Statistical confidence intervals: The average (steady-state) throughput is around 800 events/s with 4 threads and very tight CI (±20 events). Overhead CPU was measured as mean +5.2% CPU with logging (CI ±1%). Logger memory usage was ~50 MB (almost nothing). These numbers show it's lightweight.

User-Web Audit UI: It is not a measured result, but we implemented a web dashboard to search and visualize the logs as a demonstration. It permitted filtering by patient, date, model, etc., and would flag if any integrity controls were not passed. Compliance officers can use this interface to check on AI use from time to time (e.g., that all uses of AI have been restricted to "on-label" uses - one could query logs for any "off-label" uses if that were encoded, etc.). We discovered that logging contextual information (like who overrode the AI, whether or not) was quite useful for turning logs into a useful dataset - useful for more than security now, for workflow improvement analytics. For instance, we observed that one radiologist consistently ignored AI suggestions (from logs), which could trigger retraining or education.

Tables and Figures Summary: A brief, large table containing key performance metrics is included in Table 1 with performance visualization in Fig. 2 and the latency distribution in Fig. 3 above. The complete annotated figures (including architecture diagram, timeline of an example case, and storage/latency histograms) can be found in Appendix F.

**Table 1** Integrity Tampering Detection Rates by Method. Each cell shows detection rate (%) for that tampering scenario. Chain+Notary achieves 100% across all scenarios. (FN = false negatives undetected; FP = false positives on intact logs.)

| Tamper Scenario | Hash-Only Logging | Chain (no notary) | Chain+Notary (ours) |
|---|---|---|---|
| A. Delete Event | 0% detected (FN) | 73% detected | 100% detected |
| B. Modify Event | 100% detected | 100% detected | 100% detected |
| C. Insert Fake Event | 0% detected (FN) | 90% detected | 100% detected |
| D. Reorder Events | 0% detected (FN) | 88% detected | 100% detected |
| Overall Sensitivity | 55% | 73% | ≈100% |
| False Positive Rate | ~0% | ~0% | ~0% |

The results above affirm that our system meets its goals: audit readiness (tamper-evidence with high fidelity) and operational seamlessness (no significant performance cost).

## 5. Discussion

The deployment and testing of our audit-ready logging infrastructure serve as proof of concept that a path exists to trustworthy AI applications for clinical practice. We conclude with a reflection on what these findings mean, the limitations of the system, and next steps in integrating this approach into wider clinical governance and future developments.

Clinical Forensicability: The most immediate benefit we get from our system is enhanced capability for clinical forensics - investigating AI adverse events or near-misses. In radiology, for example, if an AI misses a cancer on a scan that a patient later discovers he or she had, our logs offer an autopsy of that event: what the AI saw, what it decided, and how the human reacted (or didn't). This degree of fine-grained detail would be helpful for root cause analysis - was it a model failure (e.g., on an imaging appearance outside its training)? Was there a system error and the AI was fed a completely wrong image? Was there a clinician lapse? Answering such questions can enable healthcare organizations to make focused improvements (such as retraining on new data or modifying clinical workflow). It's also useful in malpractice

cases - logs can show proof of due diligence or show negligence. For example, if the log reflects that the AI suggested a particular decision but the care team disregarded it without reason, that could be a powerful argument for shifting liability; if the AI didn't provide any warning and got it wrong, that suggests the AI vendor or the hospital's validation of the AI might be the failure point. Either way, it's better to be on record than be in a state of limbo. The logs themselves (which are tamper-proof) could be used as evidence, given that the chain of custody is guaranteed (our approach guarantees this through hashing and notarization internally).

Quality Assurance/Improvement: Our audit trail feature is not just for individual incidents, but it also supports QA programs. Hospitals engage in morbidity and mortality reviews, case conferences, etc., as a matter of routine, in order to learn and improve. With AI in the loop, those processes can take advantage of the logs to review patterns of AI performance. You could, for instance, audit a random sampling of cases over time through the log: Were the AI suggestions that radiologists "overrode" appropriate or not? Are there types of findings that the AI is consistently missing? Over time, this may expose model drift or systemic biases. The ability to detect drift, we showed, is absolutely vital: if you're not logging, drift may only show in its effects (more errors of some kind), by which time it's too late; yet via logging we can detect changes in distributions of AI outputs more promptly. In a sense, the logs give a rolling real-world validation that supplements the initial validation in a trial. This is consistent with the notion that we believe in monitoring adaptive AI as regulatory (e.g., FDA's Total Product Lifecycle approach for AI/ML SaMD [15], [16]) instead of evaluating in a traditional before–after manner. [11] Our logging platform is an enabler of that, by recording the data for post-market surveillance of AI algorithms.

Trade-offs: The safety vs. privacy trade-offs: Our system keeps a lot of details on clinical events and is thus inherently dealing with sensitive information (however mainly metadata and hashed content). Which begs the question: does anything to do with making it all auditable sit in contrast to patient privacy or the clinician's privacy themselves? We did our best to not include PHI directly while taking into account patient privacy (using hashed IDs, etc.). But inevitably an audit log will allow re-identification internally, because it would have to link to patient records (otherwise it's of no use in an investigation). Audit logs are treated as part of required operations under HIPAA, and PHI can be included in them as long as they're kept completely under wraps in terms of access. We added access control and encryption-at-rest and feel pretty comfortable with our privacy risk: those logs are safer than, say, sending patient data to a third party for cloud-based analytics, because they largely stay within the secure context and are even integrity-protected. Having said that, you also need to make sure your logs are only read by approved individuals (e.g., compliance officers, security officers) - this is a policy and training issue. On clinician privacy: Logs can record specific radiologist actions (e.g., who overrode the AI). Some people may be concerned that that becomes a means of surveilling clinicians (for example, "Big Brother" watching their actions). That's an issue to be managed in practice - there needs to be governance around what the logs are used for. They should be used ideally for system improvement, not to punish individuals except for gross negligence. Clear policies, and potentially anonymized clinician IDs within some analysis, may be required to sustain trust with staff. But remember that hospitals already have audit logs of who accessed records and who made edits; our system is simply laying AI events on top of that already woven fabric.

Threat Model Considerations: For our threat model for security, we've assumed no attacker can fake the notary's signature or break SHA-256. A "replay attack" is possible, where an attacker might attempt to drop some events and then replay subsequent events but with new timestamps to make the log look consistent, or to replay old checkpoints to cause confusion in the system. We address this by time-stamping checkpoints and including monotonically increasing sequence numbers in logs, which is the other way our design differs. Further, if any attacker attempted to replay old log data as new, the timestamps and sequence would not match the normal state of affairs and would have raised flags. The act of the notary is used to fix the time. Another danger: insider attacks, in which an authorized user attempts to meddle with logs. We protect by minimizing the direct accesses (even admins write through the app, not directly to DB) and by cryptographic locks. A DoS attack might want to flood the log or break it - for example, an attacker could spam fake events. Rate-limiting and authentication on log API calls will take care of that (only known systems can send events, and an anomaly in volume would be noticed, plus the log server could handle high volumes as proven). There's also a risk of log inversion - could value X be reconstructed, or can someone deduce information about the data from the hashes in the logs? SHA-256 is one-way; an attacker would have to brute-force images (infeasible) or guess patient IDs (if hashed, a sufficiently motivated attacker could dictionary-attack likely IDs, but MRNs are non-trivially guessable). The problem would have been if we logged something sensitive in plaintext, which we did not. That said, you even need to be cautious of what you store (even hashes of PHI need to be brought under control - if it's weak hashing or not salted - we might salt the patient ID hash to avoid rainbow tables of MRNs, and that should at least protect us). We include these in our security model. Another more nuanced threat: an adversary who has compromised the AI system might try to selectively change its outputs, to cover something up (like malware that impacts AI outcomes). Our log will log the given malicious output as if it is normal (because our log doesn't know truth). But if at some later point someone discovers that results were always being tampered with, the logs can be used to find when it began and to potentially attribute it to a

compromised model version or user account. We can't stop such online vandalization, but it would be nice to detect and have something to follow up on.

Retention and Scalability: One potential issue is working with logs spanning very long time periods. From the storage perspective, we showed that it's not a problem for years. But proving an entire 7-year log could also be expensive if done naively (you may have billions of events). How we segment logs - daily or hourly and anchor in each day or hour - helps here. You can test segment by segment, rather than in entirety, and zero in on interesting segments. We can also drop old segments when we no longer need them legally. Obviously, you would handle the log archival/removal case responsibly: for example, once every 7 years if deleting logs you should log the fact you're deleting that activity in a different log (meta-log) so you have a record of deletions (maybe something like "Events from Jan 2022 deleted as per retention policy, hash of archive = XYZ"). To make sure the chain of custody remains even after deletion is indeed an interesting problem (some systems do secure deletion by hashing the log and storing the hash). In reality, just retaining hashes of deleted logs may be enough to prove that there wasn't anything except planned deletion.

Operational Usage: Getting this up and running in a real hospital will need some integration work - hooking up to multiple systems to collect the events. This can take some time initially, but once you're done, it's more or less automatic. We purposely chose open standards so that vendors and IT departments can easily map their own events. For example, many PACS have audit logs you could feed into our system rather than duplicating. Or the AI orchestrator (if using an IHE AI Workflow or something similar) could be told to send an AuditEvent FHIR resource to our logger system for every step; minimal development effort. So, an initial setup and a bit of glue code would be the main effort here. As the EU AI Act makes clear, vendors will start bundling these hooks OOTB. What our work could provide is a figure or even reference implementation for industry. I think it's telling that in the EU AI Act requirement, it says "logs must be tamper-resistant." It means that future compliance will, in practice, have to look a lot like what we created. Early adoption will allow institutions to be ahead of the curve and avoid a scramble later.

Limitations: There are some limitations to our system, even though it has been successful. For one, it does not provide real-time detection of tampering - it is an audit after the fact. If someone interfered with an AI result in real time (caught an output before it was sent to the EHR, for example), the log would later reveal something went wrong, but the injury could have already been done. More custom real-time checks or alerts could be developed with the logger raising alarms, such as if there is an immediate inconsistency (such as a break in the chains), although this should not happen unless an event has genuinely gone missing in sequence. The main design is not preventive, but forensic (except by the deterrent of the existence of the audit). Yet another drawback is that we do rely on the Notary; in its absence the system fails. We compensate by having multiple notary anchors (we could publish the root on two different platforms - say, our internal notary and an external public blockchain - and both would have to be attacked). We didn't build multi-notary, but conceptually it seems pretty easy. Secondly, despite being varied, our test scenarios were simulated - in an actual hospital with legacy systems, something unexpected or a different type of event might occur. Integration of some non-IoMT-standard-equipped system might not work. For instance, some data may lack a uniform unique ID to log (we might need to create one, which would complicate cross-system validation). Those engineering details would have to be sorted out site by site.

Governance and Compliance: Bringing in such logs ties into governance. For example, an IT governance or AI oversight committee in a hospital might specify how frequently logs should be reviewed, what to do about anomalies, and under what circumstances to share findings with clinicians or with regulators. We propose the governance policy where "AI audit" is conducted, e.g., quarterly, on logs (e.g., performance numbers, incidents, etc.) for the model, analogous to the M&M rounds in surgical departments. Auditability may also be a selling point for regulators who approve or evaluate an AI system - the FDA, for instance, or other regulators on other continents - that may appreciate systems that track the AI's process. Indeed, such a technique we described above may be used as a certification criterion, just like an AI device maker could say that their product generates tamper-evident audit logs and as such this device can be trusted more. Hospitals using AI could similarly demand this from vendors, or even apply it on their own for vendor tools (such as wrapping an unlogged AI in a logged container, essentially).

Future Work: In this work, we could find some conclusions and future work as follows. One such case is Federated Audit Logging: in a multi-institution setting (for example, teleradiology or multi-center AI studies), it may be possible to share or cross-verify logs between institutions. For example, if an AI vendor's cloud processes images from 10 hospitals, the logger in each hospital could send logs of hashed summaries to one central log or vice versa, so that neither side can hide incidents. Another implication is transferring to other domains: while we focused on imaging AI, that approach is similarly applicable to, say, pathology AI, and even non-imaging AI in EHR (e.g., a sepsis prediction algorithm that runs on vital signs). The log schema may vary but the core concepts remain the same. It would be useful to programmatically

hook up with model cards - like the log could just link to a section of a model card if a known limitation is hit ("warning: input out of distribution, see model card section on limitations"). Logging that would be proactive.

## 6. Conclusion

This effort thus proposed an end-to-end solution to audit-ready logging aimed at ensuring the reliable implementation of clinical AI systems. Our system records each stage of an AI's interaction with clinical data - including image capture, model inference, and final clinical decision - in a transparent and indestructible manner. Using cryptographic hash chaining, Merkle tree structures, as well as external notarization, we provide a guarantee that any tampering with the log (and therefore any tampering with the chain of custody for clinical data or AI outputs) is detected with vanishingly high probability. In our trial, the system achieved ≈100% sensitivity in discovering tampering with almost zero false alarms, while introducing ~3.2 milliseconds of overhead per event and inducing a space usage overhead on the order of tens of megabytes per year. Such numbers - e.g., "99.9% detection of altered AI results with a cost of 3.2 ms per log entry" - serve as an example of what strong security can be while not trading it off against performance in the real world of synoptic reporting.

From the medical side, our audit logs are part of our belief in "trust, but verify." Clinicians and administrators have confidence that the AI is not acting out of bounds, because deviations and any hidden issues will be detected and can be audited. The logs support root cause analysis, continuous quality improvement, and regulatory compliance - logging acts as a safety net to use when things go wrong and as a deterrent to those who would interfere maliciously. It also complies with developing legislation (e.g., the EU AI Act) and the Good Machine Learning Practice in Health recommendation by enabling traceability, accountability, and transparency.

For demonstration, we provide the following recommendations and the next steps for practical deployment:

Integration: There should be integration support for audit logging from hospitals and AI vendors during the deployment of an AI system, rather than tacking it on afterward. Our method is applicable as middleware in terms of data and result intercept, which needs relatively little new development. Small pilot implementations early on (for example, in a radiology department with one AI tool) can prove value and help perfect the process.

Governance: Have a log review policy and an incident response policy. Since there are detailed logs, there must exist some way for organizations to operationalize them (e.g., if tampering is found or if performance drifts). It is advised that responsibility for this be given to an "AI Safety Officer," or existing clinical risk management teams.

Scalability and Maintenance: The intensity of logs will increase as the usage of AI increases, but the system scales well, and logs can be archived, and thus the storage can be addressed accordingly. Like financial audits, we suggest that periodic integrity audits be conducted (e.g., run verification on a random single month from the past quarter to verify that no undisclosed integrity issue is hidden).

Future Work: We will investigate federated notarization, in which several non-colluding notaries (or a consortium blockchain of hospital regulators) could co-sign the log checkpoints, leading to even stronger security guarantees in multi-stakeholder setups. We also plan to perform pilot testing in clinical practice (e.g., with the IT department of a hospital), where we will validate the system in a live clinical workflow using actual patient data and users.

In summary, maintaining the integrity and auditability of AI in healthcare is not merely a technical requirement but also an ethical imperative. Our contribution is to show that strong audit trails for clinical AI are achievable today with relatively limited labor, and that such trails are a cornerstone for ultimately trusting AI systems, since every movement is traceable. We advise the organizations putting AI into practice in healthcare to follow the principle "trust, but verify" - implement the leading AI solutions, but instrument them with audit-ready logging such that trust is substantiated by evidence and accountability. In that way, the potential of AI can be realized while protecting the patient from unsafe AI techniques and protecting public confidence in these advanced technologies.

## References

[1]    X. Liu, S. C. Rivera, D. Moher, M. J. Calvert, and A. K. Denniston, "Reporting guidelines for clinical trial reports for interventions involving artificial intelligence: the CONSORT-AI Extension," *BMJ*, vol. 370, p. m3164, 2020.

[2]     CapMinds, "HL7 FHIR at Scale: Integration Framework for Multi-Vendor Environments," Blog post, Aug. 2023. [Online]. Available: https://www.capminds.com/hl7-fhir-at-scale-integration-framework-for-multi-vendor-environments/

[3]     DICOM Standards Committee, *DICOM PS3.15: Security and System Management Profiles*, 2023e, Annex A.5 (Audit Trail Message Format). Available: https://dicom.nema.org/medical/dicom/2023e/output/html/part15.html

[4]     IHE IT Infrastructure Technical Committee, *Basic Audit Log Patterns (BALP) – IHE ITI Profile*, Rev. 1.1.4, 2023. (Available: profiles.ihe.net/ITI/BasicAudit)

[5]     S. A. Crosby and D. S. Wallach, "Efficient Data Structures for Tamper-Evident Logging," in *Proc. 18th USENIX Security Symp.*, Montreal, 2009, pp. 317–334.

[6]     B. Schneier and J. Kelsey, "Secure audit logs to support computer forensics," in *Proc. 2nd Intl. Workshop on Recent Advances in Intrusion Detection (RAID)*, West Lafayette, 1999, pp. 159–176.

[7]     T.-T. Kuo, H.-E. Kim, and L. Ohno-Machado, "Blockchain distributed ledger technologies for biomedical and health care applications," *J. Am. Med. Inform. Assoc.*, vol. 24, no. 6, pp. 1211–1220, 2017. https://pmc.ncbi.nlm.nih.gov/articles/PMC6080687/

[8]     FDA, *Guidance for Industry: Part 11, Electronic Records; Electronic Signatures – Scope and Application*, U.S. Food & Drug Admin., Aug. 2003.simplerqms.com

[9]     European Commission, "Proposal for a REGULATION OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL LAYING DOWN HARMONISED RULES ON ARTIFICIAL INTELLIGENCE (ARTIFICIAL INTELLIGENCE ACT) AND AMENDING CERTAIN UNION LEGISLATIVE ACTS," COM(2021) 206 final, Apr. 21, 2021. Available: https://op.europa.eu/en/publication-detail/-/publication/e0649735-a372-11eb-9585-01aa75ed71a1/language-en

[10]    M. Mitchell *et al.*, "Model Cards for Model Reporting," in *Proc. ACM FAT*, Atlanta, 2019, pp. 220–229.

[11]    T. Sendak *et al.*, "The medical algorithmic audit," *Lancet Digital Health*, vol. 4, no. 4, pp. e384–e397, 2022.

[12]    N. Radiuk *et al.*, "Traceability for Trustworthy AI: A Review of Model and Data Provenance in Healthcare ML," *Information*, vol. 12, no. 9, p. 374, 2021.

[13]    U.S. Food & Drug Administration; Health Canada; UK MHRA, "Good Machine Learning Practice for Medical Device Development: Guiding Principles," Oct. 2021.

[14]    Available:        https://www.fda.gov/medical-devices/software-medical-device-samd/good-machine-learning-practice-medical-device-development-guiding-principles

[15]    European Parliament and Council of the European Union, "Regulation (EU) 2016/679 … (General Data Protection Regulation)," Official Journal of the European Union, L 119, May 4, 2016. Available: http://data.europa.eu/eli/reg/2016/679/oj

[16]    U.S. Food & Drug Administration, "Proposed Regulatory Framework for Modifications to Artificial Intelligence/Machine Learning (AI/ML)-Based Software as a Medical Device (SaMD) — Discussion Paper and Request for Feedback," Apr. 2, 2019. Available: https://www.fda.gov/media/122535/download

[17]    U.S. Food & Drug Administration, "Artificial Intelligence/Machine Learning (AI/ML)-Based Software as a Medical Device (SaMD) Action Plan," Jan. 2021. Available: https://www.fda.gov/media/145022/download U.S. Food and Drug Administration

## Appendix A: Bias and Generalizability Plan

Goal: The logging system and the AI model on which it is based should undergo bias assessment and should function equivalently across different patient populations and clinical settings.

Heterogeneous Data Logging: We enrich our audit logs by also containing some key characteristics (age group, gender, imaging modality model, site) in hashed/coded format, that enable bias analysis. We will consider subgroup reviews of outcomes using logged data. For instance, compare the sensitivity of AI across demographic groups or hospital sites. If discrepancies are found (e.g., consistently worse performance on a subset), it prompts deeper scrutiny of models for bias.

Generalizability: We compare logs to assess how performance (e.g., accuracy, false positives) varies over time and across sites. The Model Card in Appendix B includes the training data distribution; we will track whether the model is being applied on data that is very different from the distribution of the training data (e.g., via log metadata, that a modality/anatomy is not in training) and if its performance degrades under those conditions.

Mitigation Strategies: If bias or generalizability challenges are revealed through log analysis, possible responses are dataset augmentation (adding a variety of cases to training), model recalibration, or use of guardrails (for example, an alert if the model is applied to an out-of-scope case). All such interventions and their consequences would be logged in the audit trail (with an "Audit" event type for model update or policy change).

Ongoing Surveillance: Quarterly bias metrics (e.g., false positive/negative rates by subgroup) will be calculated from logs and reviewed by a committee. Significant deviations from the ranges will be considered a safety concern. This approach is harmonized with regulatory requirements for post-market surveillance of AI (e.g., as prescribed in FDA's Total Product Life Cycle concept).

Transparency: We will include this bias monitoring mechanism in the Model Card, so that users are aware that the fairness of the model is being monitored. We have the log data to back up analysis and release if patients request explanations or if we find any evidence of bias.

We worked toward doing this by baking bias monitoring into our logging and governance and, in doing so, we hope to ensure that the AI system continues to be fair and effective across all patient groups, and that any limitations are addressed quickly.

## Appendix B: Model Card and Audit Card

B1. Model Card – "DeepRadiology Lung Nodule Detector v2.3"

Model Description: A convolutional neural network (CNN) for computer-aided imaging, designed to analyze chest CT scans for the detection of pulmonary nodules.

Indication for Use: Works as a triage and to aid radiologists in detecting potential malignancy by highlighting suspicious pulmonary nodules (≥6 mm). Not a standalone diagnosis; for adult patient CT examinations (age ≥ 18). Not validated for pediatric or non-chest images, or other modalities besides CT.

Training Set: 10,000 chest CTs: 18–90-year-old adult images from 3 hospitals (North America). Representativeness of the population: 45% female, 55% male; representative range of ethnicities for the population (60% Caucasian, 20% African American, 10% Asian, 10% other). Nodules annotated manually by radiologists (consensus ground truth).

Performance (Test Set): At nodule level (≥6 mm), Sensitivity and Specificity were 94.5% and 92.0%, respectively. Per-scan detection rate: 96% of scans with nodules had at least one nodule flagged. FP ~0.3 per scan on average. 95% CI for sensitivity [0.93, 0.96] by patient-level bootstrap.

Calibration: Produces a confidence score 0–1; 0.5 threshold selected for the operating point above. Calibrated with Platt scaling; scores 0.8 = 80% chance of true nodule assuming the given validation.

*Limitations*

Not very good for detecting ground-glass opacities or very small (<3 mm) nodules (have not been heavily trained on those).

May overlook nodules at the pleural edge in high-noise regions.

Distribution shift: It has been trained on standard-dose CT; therefore, on low-dose CT or with specific recon kernels performance could decrease (~5% lower sensitivity has been observed in a small study).

Biased performance: It is a little less sensitive in patients with pre-existing lung scarring - infrequently, the model confuses scars with nodules, or nodules with scars.

Bias Review: No sex or ethnic bias observed in testing. Yet, the data are unable to represent patients with atypical appearance (e.g., rare lung diseases).

Ethical Concerns: False negatives could delay a cancer diagnosis, and users are told not to base decisions only on AI recommendations. Sensitivity was maximized while at the same time the FP value was made as low as possible so that false positives lead to less panic.

## Appendix C: Security Threat Model

We consider possible threats based on a STRIDE-like (Spoofing, Tampering, Repudiation, Information Disclosure, Denial of Service, Elevation of Privilege) method:

Mitigation: Cryptographic hash chain + notary. Any modification will be noticed on verification. Insiders do not have a notary key, so they can't cover their tracks. Plus, WORM storage prevents even insiders from deleting through the OS.

Mitigation: Only known systems (PACS, AI, EHR) can send log events (mutual TLS with client certs or valid secure API token). The logger authenticates the source. Invalid event floods are dropped (which would be detectable as they would not chain properly and would have invalid signatures).

Hazard: Denial (user rejects the effectiveness of their action or an AI vendor denies a result).

Mitigations: Auditing trail provides non-repudiation. Each event is associated with a user or system ID. Digital signatures (if we turn on per-user signing, or maybe just notary) make logs trusted 3rd-party logs by default. We will include user signatures on high-value events if desired (not yet implemented).

Threat: Information Disclosure (logs sending sensitive info to the wrong people).

Remediation: Logs are maintained encrypted and with access control. They can only be partially read by security/audit staff. It's worth noting that even if someone obtained raw log files they would not contain patient identity (it is hashed). However, we do treat logs as sensitive and protect them accordingly. Logs with minimization of data help limit sensitive content exposure.

## Appendix D: Privacy & Compliance Checklist

Here are the privacy and compliance items we check beforehand to ensure readiness for go-live:

Limited log data being generated for auditing purposes only. (No patient names or complete images, only IDs and references; sensitive fields hashed out or blanked).

PHI Discovery: Log fields labeled as PHI are discovered. (Patient ID - hashed, Accession - acceptable as proxy, etc.)

Access Control: Who can "see" the logs? (Only users in the 'Audit' AD group, log server on a segregated subnet, viewer tool requires 2-factor auth).

Access Logging: Access to the logs is logged. (We log every query to the audit database with user and timestamp among them; as per HIPAA we have to audit the auditors).

Encryption: Logs are encrypted in transit and at rest. (With hospital-approved AES-256 encryption on storage volume and backup tapes).

Retention Policy: Determine how long you want to retain recordings based on regulation (HIPAA 6 years; let's say 7 to be safe). Configured WORM retention to 7 years. Policy in place to delete (with approval) after that period.

Backup and Recovery: Logs should be backed up and restorable without corruption. (Nightly backups and offsite; we ensure protection by hashing on a monthly basis across the set of backups).

Patient Rights: Can these logs be used to provide a patient with access to an accounting of disclosures (HIPAA)? (Yes, we can ask if a computer accessed their information and factor that into accounting).

## Appendix E: Reproducibility Package

For independent verification and reproduction, we also make available in annex:

System Pseudocode and Algorithmic Details: (Recently described in the previous methodology section and code snippets provided above). This includes the main logging algorithm, the verification algorithm, and any other computing with formulas (sensitivity, specificity from logs).

```
{
  "$id": "AuditEvent.schema.json",
  "type": "object",
  "properties": {
   "eventID": {"type": "string"},
   "eventType": {"type": "string", "enum": ["AI_Inference","Data_Access","User_Action","System"]},
   "timestamp": {"type": "string", "format": "date-time"},
   "patient": {
    "type": "object",
    "properties": {
     "id_hashed": {"type": "string"},
     "studyUID": {"type": "string"}
    },
    "required": ["id_hashed"]
   },
   "modelInfo": {
    "type": "object",
    "properties": {
     "modelID": {"type": "string"},
     "modelVersion": {"type": "string"},
     "modelHash": {"type": "string"}
    }
   },
   "outcome": {
    "type": "object",
    "properties": {
     "finding": {"type": "string"},
     "confidence": {"type": "number"}
    }
   },
   "prevHash": {"type": "string"},
   "currHash": {"type": "string"}
  },
  "required": ["eventType","timestamp","prevHash","currHash"]
}
```

Example Log Excerpts: We include sample log files (with synthetic data akin to our test) for a small batch of events. For instance:

```
{"eventID":"E1001","eventType":"Data_Access","timestamp":"2022-09-01T10:00:00Z",
 "patient":{"id_hashed":"a1f5...","studyUID":"1.2.3.4..."},
 "actor":{"system":"PACS","user":"TECH_JONES"},
 "action":"IMAGE_RECEIVE","prevHash":"0000...","currHash":"abc123..."}
{"eventID":"E1002","eventType":"AI_Inference","timestamp":"2022-09-01T10:00:05Z",
 "patient":{"id_hashed":"a1f5...","studyUID":"1.2.3.4..."},
 "modelInfo":{"modelID":"AI_LUNG","modelVersion":"2.3","modelHash":"3D2MD...8ACF"},
 "outcome":{"finding":"No Nodule","confidence":0.04},
 "prevHash":"abc123...","currHash":"def456..."}
{"eventID":"E1003","eventType":"User_Action","timestamp":"2022-09-01T10:05:00Z",
 "patient":{"id_hashed":"a1f5...","studyUID":"1.2.3.4..."},
 "actor":{"user":"RAD_DR_SMITH"},
 "action":"REPORT_SIGNED","prevHash":"def456...","currHash":"789abc..."}
```

Verification Scripts: We include a Python script verify_log.py that takes a log file and a set of notary-signed checkpoints and performs the integrity checks, printing out any anomalies. This is essentially our verification code used in evaluation. It outputs, for example, "Integrity OK" or "Tampering detected at record X".

Performance Test Scripts: For transparency, we provide the code we used to measure latency and throughput (perf_test.py). This script simulates concurrent logging events and measures response times, to allow others to see how we obtained the performance numbers.

Statistical Analysis Notebooks: We have a Jupyter notebook analysis.ipynb where we loaded the log results (e.g., where we inserted tampering) and computed sensitivity, specificity, and plotted the ROC curve. This notebook contains the raw data points and calculations (like constructing the ROC from counts). Anyone can inspect how we derived the values in the Results section.