



White Paper

AI-ENHANCED LONGITUDINAL PATIENT PROFILES: Connecting Episodes of Care Across Settings

WWW.DEEPKNIT.AI

TABLE OF CONTENTS

1. Executive Summary	3
2. The Fragmentation Crisis: An Operational and Clinical Liability	4
3. The Economic Impact of Disconnected Data	5
4. The Human Cost: Patient Safety	5
5. The Evolution of Identity: From Probabilistic to Referential Matching	6
- Limitations of Legacy Methodologies	
- The Third Generation: Referential Matching	
- AI and Machine Learning: Automating Stewardship	
6. Privacy-Preserving Record Linkage (PPRL): The Key to Research	8
- The Mechanics of Bloom Filters	
- Empowering Real-World Evidence	
7. Clinical Application: Chronic Care and Coordination	9
- Case Study I	
- Case Study II	
- The Role of Generative AI in Synthesis	
8. Legal, Ethical, and Regulatory Frameworks	10
- TEFCAs and the QHIN Mandate	
- Liability and the “Moral Crumple Zone”	
- Addressing Bias	
9. Conclusion	11
10. References	12





EXECUTIVE SUMMARY

The United States healthcare infrastructure is currently generating data at an exponential rate (50 petabytes/hospital annually), yet the industry's ability to synthesize this information into a coherent narrative for the individual patient remains critically underdeveloped.

A patient's medical history is rarely a single, cohesive document; rather, it is a scattered archipelago of isolated data points—across clinics, hospitals, labs, and telehealth, creating data silos that hamper care. This fragmentation poses a profound threat to patient safety, a primary driver of administrative waste, and the single greatest technical barrier to value-based care.

This whitepaper explores the transformative role of Artificial Intelligence (AI) and Machine Learning (ML) in resolving this crisis of identity. We examine the industry's necessary shift from traditional probabilistic matching to AI-driven Referential Matching, the deployment of Privacy-Preserving Record Linkage (PPRL) for research and population health, and the legal frameworks required to govern these technologies. By leveraging AI to create a "Single Source of Truth," healthcare organizations can move beyond simple identity verification to true care continuity, reducing the liability risks for clinicians and ensuring that the longitudinal patient profile becomes a reliable instrument of care.



THE FRAGMENTATION CRISIS: AN OPERATIONAL AND CLINICAL LIABILITY

The concept of the “longitudinal patient profile” represents the ideal state of health informatics: a chronologically ordered, comprehensive dataset capturing every interaction a patient has with the health system, from birth to death. However, the reality of the U.S. healthcare system is defined by siloes. Data resides in proprietary Electronic Health Records (EHRs), laboratory information systems, payer claims databases, and increasingly, patient-generated health data (PGHD).

The core impediment to unifying this data is Identity Resolution—the process of determining that “John Smith” at an urgent care clinic, “J. Smith” at a radiology center, and “Jonathan Smith” in a payer database are the same biological entity. Without accurate identity resolution, the longitudinal record is impossible to construct.

Micky Tripathi, the National Coordinator for Health Information Technology, has analogized the current state to a cellular network where carriers are not connected—users can make calls within their network but are isolated from the broader system. The Trusted Exchange Framework and Common Agreement (TEFCA) aims to bridge these networks, but its success relies entirely on the technical ability to match identities accurately across disparate systems.



THE ECONOMIC IMPACT OF DISCONNECTED DATA

The costs of failing to link patient records are actually mind numbing. Duplicate records are often described as the “dark matter” of healthcare - unseen but exerting a massive gravitational pull on efficiency.

- Operational Waste:** It costs a healthcare organization approximately \$1,000 per pair to reconcile duplicate records manually. For large systems, this creates millions of dollars in administrative overhead.
- Revenue Leakage:** Roughly 30% to 35% of all denied claims in the U.S. are attributed to inaccurate patient identification or incomplete information. The Ponemon Institute estimates this costs the average healthcare system \$1.2 million annually.
- Redundant Care:** When a provider cannot access a patient's prior history, they often repeat expensive diagnostic tests. This not only incurs unnecessary costs but also exposes patients to avoidable risks, such as radiation from repeated CT scans.



THE HUMAN COST: PATIENT SAFETY

The clinical consequences of fragmentation are severe and actionable under malpractice law. A survey by the College of Healthcare Information Management Executives (CHIME) revealed that one in five hospital CIOs reported patient harm due to mismatches in the previous year.

Duplicate records are associated with a 1.44 times higher odds of missing a critical laboratory result. If a physician treats a patient under one Medical Record Number (MRN) while a critical lab value or allergy is stored under a duplicate MRN, the result can be fatal. In fact, duplicate records account for nearly 2,000 preventable deaths each year in the United States.

Scenario	Statistic	Implication
Denied Claims	35%	Direct revenue loss due to identification errors.
Match Rates (Internal)	~80%	1 in 5 patients may have a fragmented record within a single facility.
Match Rates (External)	~50%	Inter-organizational data exchange fails half the time without advanced matching.
Harm Frequency	20%	Percentage of CIOs reporting patient harm due to mismatches.
Mortality	~2,000/yr	Annual preventable deaths attributed to patient matching errors.



THE EVOLUTION OF IDENTITY: FROM PROBABILISTIC TO REFERENTIAL MATCHING

To solve fragmentation, the industry has evolved through several generations of matching methodologies. We are currently witnessing a transformative shift from internal, rule-based systems to external, data-driven AI architectures.

- Limitations of Legacy Methodologies

For decades, the industry relied on Probabilistic Matching, often based on the Fellegi-Sunter model. This approach assigns weights to different fields based on the likelihood that they indicate a match (e.g. matching a rare name like “Zebulon” carries a higher weight than matching “John”). While effective in restricted environments, probabilistic matching struggles with the dirty, static nature of real-world healthcare data. It creates a massive gray area of potential matches that requires human oversight, and its performance drops drastically when data crosses organizational boundaries, often yielding match rates as low as 50-60%.

- The Third Generation: Referential Matching

The most significant leap in identity resolution is Referential Matching. This approach acknowledges that hospital data is inherently messy. Instead of comparing two inconsistent hospital records against each other, referential matching compares both records against a massive, curated external database—a “Reference Database” or “Referential Knowledge Base”.

Vendors like Verato and Experian aggregate data from diverse, non-healthcare sources—credit headers, utility records, public records, and USPS data—to build a comprehensive “life history” of individuals. This allows the system to bridge gaps that probabilistic algorithms cannot. For example, if a patient record at a lab lists “Mary Smith” at an old address, and a hospital record lists “Mary Jones” at a new address, a probabilistic engine would likely fail to link them.

However, the referential database contains the historical knowledge that Mary Smith changed her name to Jones and moved in 2022, enabling a confident match. Studies comparing these methodologies have shown that referential matching improves sensitivity from approximately 64% (probabilistic) to over 93–95%, drastically reducing the creation of duplicate records.

“Referential matching acts as an answer key. Instead of trying to guess if two records are the same based on limited data, you check them against a complete, historical profile of the person’s identity data spanning decades.” Industry Analysis of Verato Methodology.

- AI and Machine Learning: Automating Stewardship

Artificial Intelligence enhances these methods by automating the complex “human judgment” required to resolve edge cases. AI models, particularly those using Random Forest classifiers and Deep Learning embeddings, can analyze the non-linear relationships between data points. These models learn subtle patterns; for instance, that a mismatch in a middle initial is less significant if the SSN and address are identical, but highly significant if the address is a generic hospital address.

This capability enables “Auto-Stewardship.” In traditional systems, the “gray area” of potential duplicates accumulates in queues that require manual review by Health Information Management (HIM) staff. AI models trained on millions of previous human decisions can resolve the vast majority of these ambiguous cases automatically. This reduces the administrative burden on HIM departments by orders of magnitude, freeing staff to focus only on the most complex forensic identity issues.



PRIVACY-PRESERVING RECORD LINKAGE (PPRL): THE KEY TO RESEARCH

As healthcare moves toward population health and large-scale research, the need to link records without exposing Personally Identifiable Information (PII) is paramount. Privacy-Preserving Record Linkage (PPRL) is the technical architecture enabling this secure synthesis.

- The Mechanics of Bloom Filters

PPRL allows two organizations, such as a genomic research center and a health system, to determine shared patients without sharing patient lists. It uses hashing and Bloom Filters: while standard hashing (e.g., SHA-256) is brittle because a single typo changes the entire hash, Bloom filters split PII into “q-grams” (character sub-sequences) and map them to a binary vector. Comparing these vectors helps identify matches despite slight data variations, while the data remains encrypted.

- Empowering Real-World Evidence

This technology is critical for generating Real-World Evidence (RWE). For example, Flatiron Health utilizes longitudinal datasets derived from PPRL to study rare hematologic malignancy subgroups. By linking clinical EHR data with genomic data and claims history, researchers can track patient trajectories across years, identifying outcomes for treatments in populations that are historically underrepresented in clinical trials. Similarly, during the COVID-19 pandemic, PPRL allowed public health agencies to link vaccination registries with hospitalization records to determine vaccine efficacy without compromising patient trust or violating HIPAA.



CLINICAL APPLICATION: CHRONIC CARE AND COORDINATION

The true value of the longitudinal profile is realized in the management of chronic conditions, where the patient's journey spans multiple years and dozens of providers.



Case Study I

A leading integrated delivery network healthcare in New York struggled to exchange data with the University of Utah. Despite using the same EHR vendor, match rates were as low as 10% due to data inconsistencies. After implementing a referential matching layer, the match rate rose to 95%, enabling plug-and-play interoperability and allowing queries to locate patient records across disparate systems without manual intervention. Additionally, the AI-driven "Auto-Steward" increased data stewardship productivity by 174%, resolving thousands of potential duplicates overnight and ensuring clinicians had access to complete histories for complex patients.



Case Study II

One of the largest providers in NY faced a backlog of 200,000 potential duplicate records after rapid expansion and acquisitions, risking revenue cycle operations and patient care. Using advanced matching algorithms and a dedicated data integrity team, the backlog was cleared. The resulting clean Master Patient Index (MPI) became the foundation for digital transformation, enabling a reliable "patient 360" view for billing and clinical decision support.



The Role of Generative AI in Synthesis

Once the longitudinal record is assembled, the challenge shifts from access to cognition. A physician cannot read thousands of compiled documents in a fifteen-minute visit. Large Language Models (LLMs) are now being deployed to summarize these longitudinal profiles. For example, AI models can generate "problem-based summaries" that extract and synthesize only the relevant history for a specific condition, such as heart failure, while filtering out unrelated data.

This reduces cognitive load and ensures historical data informs current clinical decisions. Benchmarks like Stanford's EHRSHOT are setting standards for how models predict patient trajectories from longitudinal data, moving AI from backend infrastructure to frontline clinical support.



LEGAL, ETHICAL, AND REGULATORY FRAMEWORKS

The deployment of these technologies operates within a complex web of regulation and liability.

- TEFCA and the QHIN Mandate

The Trusted Exchange Framework and Common Agreement (TEFCA) establishes a “network of networks” for nationwide connectivity through Qualified Health Information Networks (QHINs), which route data across the country. Its QHIN Technical Framework (QTF) requires strong Patient Identity Resolution, ensuring a QHIN can determine with high probability that a patient query matches a record at a distant organization. This federal mandate is driving referential matching as the standard for interoperability.

- Liability and the “Moral Crumple Zone”

As AI increasingly builds patient profiles, a key legal question emerges: who is responsible when it's wrong? Under current negligence law, clinicians are often liable as the learned intermediary, even if the error originates in the software—creating a “moral crumple zone” where humans absorb the machine's mistakes. Legal scholars are proposing “risk pooling” or “shared responsibility” models, where AI developers and vendors share liability for matching errors, encouraging them to prioritize safety and accuracy over speed.

- Addressing Bias

AI models are only as good as their training data. Traditional algorithms often struggle with multi-part Hispanic surnames or Asian naming conventions, leading to higher error rates for these populations. It is therefore an ethical and clinical imperative to audit modern AI matching systems for demographic bias, ensuring longitudinal records remain accurate and accessible for all patient groups.



CONCLUSION

The creation of AI-enhanced longitudinal patient profiles is not merely a technical upgrade; it is a fundamental restructuring of how healthcare understands the individual. By moving from fragmented episodes to a cohesive narrative, we empower clinicians with the context required to save lives. The technology exists: Referential Matching provides the accuracy, PPRL provides the privacy, and Generative AI provides the insight.

The barriers remaining are largely cultural and political. Achieving the vision of a truly connected healthcare system requires a concerted effort to embrace these technologies, enforce the standards of TEFCA, and accept that in a digital age, accurate identity is the most critical medical instrument of all. As we look to the future, the “Longitudinal Patient Record” will cease to be a static document and will become a dynamic, AI-driven model of the patient—a digital twin that travels with them, ensuring that no matter where they seek care, they are known, understood, and treated with the full weight of their medical history.



REFERENCES

1. **The role of artificial intelligence for the application of integrating electronic health records and patient-generated data in clinical decision support** – [PubMed Central](#)
2. **Longitudinal care records: How an EMPI can help make these a reality** – [Rhapsody](#)
3. **A TEFCA Primer with Micky Tripathi** – [This Week Health](#)
4. **2024 Costs of Caring** – [AHA - American Hospital Association](#)
5. **Identity Resolution and Data Quality Algorithms for Person Indexing** – [Oracle Health Sciences](#)
6. **Enhanced Patient Matching Is Critical to Achieving Full Promise of Digital Health Records** – [The PEW Charitable Trusts](#)
7. **Referential Matching Could Be the Answer to Better Patient and Member Identification** – [LexisNexis Risk Solutions](#)
8. **Advancing Responsible Healthcare AI with Longitudinal EHR Datasets** – [Stanford University: Human-centered Artificial Intelligence](#)
9. **Verifiable Summarization of Electronic Health Records Using Large Language Models to Support Chart Review** – [MedRxiv](#)
10. **Patient Identification Techniques – Approaches, Implications, and Findings** – [PubMed Central](#)