# Learning-Based Semantic Alignment in Cross-Institution Health Data Exchange

Gayathri Surianarayanan
Independent Researcher
Virginia, USA
https://orcid.org/0009-0001-7987-5090

*Abstract*—Healthcare ecosystems today are built upon the exchange of health data between multiple types of organizations, including provider organizations, payers, laboratories, public health agencies, and government programs. While syntactical standards, such as HL7 v2, FHIR, X12, and CDA have been widely adopted to enable the exchange of data electronically, the lack of consistent semantics across institutions continues to create barriers to the correct interpretation of data exchanged electronically, automated workflows and analytics. The variability in schema design, terminology use, contextual meaning, and workflow specific customization creates ambiguity in data exchanged electronically which cannot be eliminated by using rule based mapping techniques or static ontologies. In this paper we present a novel, learning-based semantic alignment framework for cross-institutional health data exchange which utilizes representation learning, contextual embeddings and probabilistic alignment scoring to align schematically disparate healthcare datasets. Our proposed approach consists of three main components; schema-level feature learning to extract common features from the source and target schemas, a context aware similarity model that measures the similarity between the two schemas, and an explainability driven validation process to ensure that the alignment decisions made are trusted and suitable for regulated environments. We also present a multi-layer architecture to facilitate both semantic discovery, alignment confidence estimation, human-in-the-loop review and regulatory auditability. We evaluate our proposed approach on synthetic but realistic healthcare schemas developed from electronic health record (EHR) data, claims data and enrollment system data, and compare our approach with two baseline approaches; one that uses rules to align the two schematics and another that uses an ontology to align the two schematics. Finally, our proposed approach includes governance controls that meet all applicable requirements of HIPAA, CMS, and NIST AI Risk Management Framework (RMF).

*Index Terms*—Semantic Interoperability, Health Data Exchange, Machine Learning, Schema Alignment, Healthcare Interoperability, FHIR, Claims Data, Data Governance, Explainable AI, Regulatory Compliance, HIPAA, NIST AI RMF

## I. INTRODUCTION

Data sharing is a core function of healthcare data exchanges. Care Coordination, Claims Adjudication, Eligibility Determination, Quality Measurement, Population Health Analytics, and Regulatory Reporting are all examples of how data sharing is used in healthcare. Standards for data sharing have improved significantly over the last decade with the advent of HL7, FHIR, X12, CDA, etc. Data exchange at scale was possible because these standards defined both the format of the data and how it would be transported between systems. However, syntactic interchangeability alone does not ensure semantic interchangeability. In many cases, organizations will take the standard schema and either modify the standard by changing the semantics of specific fields, add their own extension, or interpret the standard in a way that fits their organization's workflow. When different organizations use the same terms but have different meanings, they create inconsistent analytics, incorrect policy interpretation, and ultimately failure of downstream automation processes. The traditional methods of achieving semantic alignment are primarily manual mappings, expert driven ontology alignment, or deterministic transformation pipeline. Although successful in small, and relatively static environments, the scalability, flexibility to accommodate changes in schemas, and the ability to model ambiguities in semantic relationship are major challenges in large-scale cross-organization data exchanges. As the number of data sources increases, so does the cost of manual maintenance, and the accuracy of the mappings decreases as the schemas evolve and the workflows change. In order to address these challenges, this paper presents a learning-based semantic alignment framework for cross-institutional healthcare data exchanges. This includes schema representation learning, contextual similarity modeling, and alignment confidence estimation integrated into an architecture which provides governance awareness for regulated healthcare environments. A diagrammatic representation of the proposed semantic alignment architecture, along with an overview of the proposed framework as part of cross-institutional data exchange workflows, is presented in Figure 1. The primary research question addressed in this paper is: How can learning-based semantic alignment techniques be developed to provide reliable reconciliation of disparate healthcare schemas across institutions, while maintaining transparency, accuracy, and regulatory compliance? To answer this research question, we propose a learning-based semantic alignment framework for cross-institutional health data exchanges.

## II. BACKGROUND AND MOTIVATION

### A. Semantic Interoperability Barriers in Healthcare

The primary factor in healthcare's heterogeneity is the fact that providers, payers, labs, and government entities function as separate entities with different operational (e.g., staffing models), regulatory (e.g., HIPAA compliance) and financial
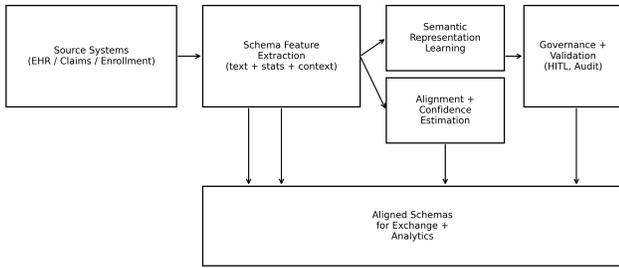
Fig. 1. Learning-based semantic alignment architecture for cross-institution health data exchange.

(e.g., reimbursement structures) constraints and therefore create disparate data representations for identical standards. Some of the most common barriers to semantic interoperability include:

- Adaptation of schemas: Each institution customizes or constrains the use of the standard schemas in order to fit their specific operational needs.
- Terminology variance: Various coding systems or localized interpretations of standardized codes may be employed by institutions.
- Dependent on Context: The semantics of a field will vary based upon its position within a hierarchy of resources, the presence of other attributes, or institutional policies and business rules.
- Schema Evolution: Frequent changes occur due to new regulatory requirements, new benefit designs and changing clinical practices.

As a result, the issues of alignment between various data representations can be very difficult to identify using static rules; specifically where there are slight semantic differences versus significant syntactical differences. In addition, the architecture illustrated in Figure 1 demonstrates the process of transforming various, heterogeneous healthcare schemas into one unified representation of those schemas via multi-modal feature extraction and learning-based alignment. In doing so, the architecture addresses several limitations associated with traditional rule-based and ontology-only methods of representing semantic interoperability, such as schema drift and variability of institutional settings.

### B. Limitations of Rule-Based and Ontology-Driven Alignment

Generally, rule-based alignment is limited by its dependency upon the string matching of rules (as in template matching), predetermined transformations (as in a set of predetermined mapping templates) and/or pre-constructed mapping tables (as in "hand-coded" mapping tables). Ontology-based alignment utilizes a controlled vocabulary that captures hierarchically related terms, and uses these hierarchical relationships to determine which terms have similar meanings. Both rule-based and ontology-based alignment have their own limitations:

1) Limited ability to adapt to changing schemas: The alignments will generally require manual maintenance and modification as the schemas change over time.
2) Poor generalization: Rule-based and ontology-based alignment methods often perform poorly when the schema includes new terms and/or other unseen naming conventions.
3) Lack of an estimate of confidence: Since rule-based and ontology-based alignment methods always generate deterministic matches, there is no measure of how confident one should be in those matches, thus increasing the risk of using automated tools for alignment.
4) Scalability limitations: The amount of manual effort required to maintain alignments increases exponentially with the number of data sources involved.

### C. Motivation for Learning-Based Semantic Alignment

To derive the meaning of the semantic relationship between schema elements from the perspective of the learner, Learning-based Semantic Alignment uses both descriptive, quantitative, and contextual features (i.e., Textual Descriptions, Data Distributions, Relational Context and Usage Patterns) of Schema Elements as input for the learning process. In contrast to static rule-based systems, Machine Learning Systems can:

- Capture the implicit semantic relationships that exist between heterogeneous Schemas.
- Learn from changes in the schema over time by retraining or incrementally learning from new data.
- Assign a level of certainty to the degree of semantic alignment between Schemas to assist with Governance and Human Oversight.
- Enable Explainable Mechanisms to provide auditing and compliance with Regulatory Review.

The need for such features is especially critical in Healthcare environments; as incorrect semantic alignment can lead to financial losses, regulatory violations and patient harm. The purpose of this paper is to establish a structured and compliant learning-based semantic alignment framework, building on the previous work motivating these needs.

### D. Contributions of This Paper

The following are the significant contributions to this study.

1) Presents a multi-layer learning based on semantic alignment of cross-institutional healthcare data exchange architecture.
2) Introduces schema representation that is aware of the context that combines textual, statistical, and relational characteristics.
3) Describes alignment confidence and explainability mechanisms that can be applied in the regulatory environment.
4) Evaluate the proposed method by applying realistic healthcare schemas, and compare the results obtained from the proposed method and baseline methods.
5) Alignment with HIPAA, CMS, and NIST AI Risk Management Requirements has been demonstrated.

## III. RELATED WORK

Semantic interoperability and schema alignment have remained major hurdles for many years in healthcare information technology. The previous research studies various types of interoperability including: ontology based interoperability, rule based schema matching and more recent machine learning based schema alignment. However, there have been few approaches with respect to the application to the large scale, dynamic, regulated and multi-institutional cross-institutional healthcare environments.

### A. Ontology Based Semantic Interoperability

Ontology based approaches use standardized vocabulary and structured hierarchical representations of knowledge (e.g., SNOMED CT, ICD, LOINC, RxNorm) to bring together the health care data and provide an integrated representation of data from different systems. The ontology based approaches use standardized vocabularies and structured hierarchical representations of knowledge (e.g., SNOMED CT, ICD, LOINC, RxNorm). They attempt to map the schema elements to the common semantic concepts defined in the ontology, and therefore they attempt to align the schema elements based on the common semantic concepts defined in the ontology. The ontology based approaches are very successful for the clinical terminology normalization. However, ontology based approaches also have some significant limitations:

- There are many cases where the ontology is incomplete, or the ontology is too slow to represent new operational data models.
- The administrative and financial data (claims, enrollment, eligibility) are poorly represented in the ontology.
- It is very difficult to represent the local extension and institutional semantics in the formal way.
- The ontology based alignment do not consider the usage patterns and context dependent relationships between the fields.

Therefore, it can be said that the ontology based alignment is not sufficient to achieve the full semantic interoperability among the different healthcare organizations.

### B. Rule Based and Heuristic Schema Matching

The rule based schema alignment techniques use the deterministic logic (string similarity, data type matching, naming convention, etc.) and the manually created transformation rules to align the schema elements. Due to their transparency and predictability, the rule based systems are widely used in the enterprise integration platforms and data warehouse applications. However, as indicated by the numerous empirical studies, the rule based systems degrade quickly under the pressure of the schema drift and scalability. Therefore, it is impractical to maintain the accurate rule set for the large number of the institutions participating in the health care exchanges; each institution may have dozens of the schema variations; and the regulations for the health care exchanges may change frequently. Furthermore, since the rule based systems are deterministic, they cannot express the ambiguity

or partial semantic match; and this makes the potential of silent misalignment of the downstream analytics even higher.

### C. Machine Learning for Schema and Semantic Alignment

Recently, researchers proposed the machine learning based schema matching using the supervised, semi-supervised and unsupervised techniques. Typically, the machine learning based schema matching uses the textual metadata, instance level data, and graph-based features to create the vector representation of the schema elements and then computes the similarity score between the pair of the schema elements. Machine learning based approach provides many advantages over the traditional approaches:

- The ability to find the generalized similarity between the schema elements
- The adaptation to the schema changes
- The integration of the multiple semantic signals
- The probabilistic confidence in the alignment

Although the above mentioned advantages of the machine learning based approach are clear, most of the prior works on machine learning based schema matching are focused on the generic enterprise schemas and open datasets, and therefore, the machine learning based approach did not address the unique characteristics of the regulated, governed, and audited environment of the health care data exchange.

### D. Explainability and Trust in Data Alignment

Explainability has become one of the required properties for the AI systems working in the high-risk domain. In the context of semantic alignment, the stakeholders need to know why two schema elements are equivalent, partially equivalent, or completely non-equivalent. While the majority of the existing explainability approaches focus on the predictive model, the existing approaches do not address the structural alignment decisions. The token level attention or feature attribution is not sufficient to provide the evidence for the stakeholders to link the alignment decisions to the schema context, the distribution of the data, and the regulatory constraints. This deficiency motivated the researchers to develop the explainable alignment mechanism and integrate the explainable alignment mechanism in the machine learning based framework.

### E. The Existing Research Gaps

According to the literature review, there are several significant gaps in the existing research:

1) The limited research focus on the healthcare-specific administrative and claims data.
2) The limited handling of the schema drifts among the institutions.
3) The limited modeling of the alignment confidence and uncertainty.
4) The limited incorporation of the explainability and governance controls in the alignment process.
5) The lack of regulatory aware validation methodologies.

This paper will close the above-mentioned gaps by developing the machine learning based semantic alignment framework

specifically designed for the cross-institutional healthcare data exchange.

## IV. Proposed Framework

### A. The Problem Statement

Representing the elements of the schema of a source health care institution ($S$) and the elements of a target institution ($T$), we may define the semantic alignment problem as finding a mapping function $f : S \rightarrow T \cup \{\emptyset\}$, such that each element of the source can be mapped to zero or more semantically similar or related elements of the target or to an empty set. Due to many factors in cross-institutional health care settings, the semantic alignment problem has been particularly difficult:

- Partial overlap between schemas
- Contextualized definitions of terms
- One-to-many, many-to-one relationships
- Evolving schemas and policies
- Constraints to use data based on regulatory requirements

Therefore, the goal is to achieve semantic consistency which maintains meaning during clinical, operational and financial workflows and not just a syntactical correspondence.

### B. Design Objectives

Based on the objectives stated above, the proposed framework will meet the following design objectives:

- Semantics Validity: Maintain the true meaning of schema elements from different organizations.
- Adaptation: Adapt to changes in the structure of the schema without requiring significant manual efforts to do so.
- Interpretability: Allow users to understand why certain decisions were made.
- Awareness of Confidence: Measure uncertainty and allow human intervention when needed.
- Scalability: Allow operation over a large number of schemata and institutions.
- Compliance with Regulatory Requirements: Provide mechanisms to support audit trails and compliance requirements.

### C. Functional Requirements

In order to accomplish its design objectives, the framework should include the following functional requirements:

- Schema Representations: Multi-modal representations including both textual and statistical and relational aspects of schemata.
- Probability-based Similarity Scoring: In order to make probabilistic similarity assessments instead of making binary decisions regarding semantic similarity.
- Alignment Acceptance Thresholds and Escalation: Mechanisms to determine whether alignments should be accepted or escalated to humans.
- Validation by Humans of Low-Confidence Alignments: Mechanisms to validate alignments that have low confidence associated with them.

- Alignments Artifacts: Traceable artifacts of alignments that are suitable for auditing purposes.

### D. Non-Functional Requirements

In addition to satisfying the functional requirements, the framework must also satisfy the typical non-functional requirements of health care systems:

- Preservation of Privacy: The alignment must occur at the schema level and cannot expose Protected Health Information (PHI).
- Efficiency of Performance: The alignment latency must support batch and near real-time workflows.
- Extensibility: Ability to add new institutions and schema versions incrementally.

### E. Evaluation Criteria for Alignments

In order to establish a high degree of trustworthiness and reliability in the framework's ability to perform alignments, the output of alignments must be evaluated against the following criteria:

- Precision, Recall, and F1-Score for Semantic Mappings
- Stability under Schema Drift
- Accuracy of Downstream Analytics and Reporting Results
- Completeness of Explainability Coverage and Audit Trails

These criteria will guide the evaluation method that is described in the next section of this paper.

The proposed framework for learning-based semantic alignment for cross-institutional health data exchange enables reconciling different schemas by generating a representation of schema elements' semantics, assessing alignment confidence, and ensuring explainability and governance.

### F. Framework Overview

The proposed framework has four main layers:

1) Schema Feature Extraction Layer
2) Semantic Representation Learning Layer
3) Alignment and Confidence Estimation Layer
4) Governance and Validation Layer

Collectively, these layers enable an alignment pipeline to transform heterogeneous schemas into aligned, semantically coherent representations enabling interoperable analytics and exchange. The architecture is illustrated in Figure 1, which shows how this proposed framework fits within cross-institutional data exchange workflows. The framework only operates at the schema and metadata level, and therefore does not need access to record-level Protected Health Information (PHI), and therefore can support privacy-preserving alignment within regulated environments.

### G. Schema Feature Extraction Layer

The first layer is responsible for extracting multi-modal features from schema elements coming from different health care systems including Electronic Health Records (EHRs), claims platforms, enrollment systems, and public health data

repositories. Each schema element $e$ is extracted using the following three categories of features:

1) Textual Features:
   - Field name(s) and alias(es)
   - Description/documentation text of field(s)
   - Standard identifier(s) (FHIR path(s), X12 segment(s))
2) Statistical Features:
   - Data type(s)/format(s)
   - Value distribution(s) (range(s), cardinality, entropy)
   - Null rate(s) and sparsity indicator(s)
3) Contextual Features:
   - Relationship of parent/child among schema hierarchy elements
   - Co-occurrence of neighboring field(s)
   - Workflow specific grouping(s) (e.g., eligibility, claims, billing)

The features described above capture additional semantic information beyond what is captured by simply comparing strings. Contextual features are particularly useful in health care because much of the meaning associated with structure comes from positional relationships rather than naming conventions.

### H. Semantic Representation Learning

Using a common embedding model, the extracted features are mapped into a set of dense vector representations (i.e., embeddings). Each schema element is represented by a single composite vector representing the semantic properties of the element:

$$v_e = \alpha v_{\text{text}} + \beta v_{\text{stat}} + \gamma v_{\text{context}}$$

Where $\alpha$, $\beta$, and $\gamma$ represent tunable weights indicating the relative influence of each of the feature modalities. Textual embeddings capture similarity based on the linguistic characteristics of the feature values; statistical embeddings capture similarity based on the behavior of the feature values; and contextual embeddings capture similarity based on the relationship of the feature values to other elements in the schema. The ability to combine all three forms of similarity enables the model to find equivalences in elements that have similar meanings but have very different syntax, i.e., institutional-specific identifiers or local benefit attribute definitions. The embedding models used will depend upon the availability of data, the nature of the available data, and the applicable governance constraints. These include training using contrastive learning, weak supervision from a known mapping, and unsupervised similarity metrics.

### I. Alignment and Similarity Modeling

Alignment is modeled as a similarity estimation problem between the source and target schema elements. A similarity score is estimated as follows:

$$A(s_i, t_j) = \text{sim}(v_{s_i}, v_{t_j})$$

where $v_{s_i}$ and $v_{t_j}$ denote the embedding vectors of the source element $s_i$ and the target element $t_j$, respectively. $\text{sim}(\cdot)$ denotes the cosine similarity metric or a learned similarity function. The proposed framework supports:

- one-to-one alignments between identical fields
- one-to-many alignments when semantics are decomposed
- many-to-one alignments when semantics are aggregated

Similarity scores are used to rank candidate alignments and only the top-ranked candidates, above predefined thresholds, are accepted without human intervention.

### J. Alignment Confidence Estimation

Unlike many deterministic mapping frameworks, the proposed framework models alignment confidence. Confidence in alignment is estimated from:

- Magnitude of similarity score(s)
- Agreement of feature values across modalities
- Stability of the alignment over perturbation(s) (e.g., renaming, reordering)

Estimates of alignment confidence have several uses:

- Triggering human evaluation for uncertain mappings
- Supporting risk assessments for subsequent automation
- Providing evidence for auditing and governance processes

Alignments with low confidence are manually validated by data steward(s); high-confidence alignments may be auto-approved subject to pre-defined governance policies.

### K. Explainability and Interpretability

Explainability is integrated into the alignment process and for each alignment decision an explanation artifact providing:

- Contribution of textual, statistical, and contextual feature values
- Dominant similarity driver(s)
- Any conflicting signals, if present

These explanation artifacts provide stakeholders with insight into why a pair of schema elements were aligned, and they permit stakeholders to assess whether the alignment decision was semantically correct in their particular regulatory or operational environment.

### L. Governance and Validation Layer

The final layer integrates alignment output with organizational governance processes. These include:

- Version control of schema alignments
- Audit log entries recording alignment decisions and confidence levels
- Manual review approval workflow(s)
- Traceability to relevant regulatory requirements and data standards

By integrating governance controls into the alignment framework, the proposed framework enables continuous interoperability, while remaining compliant with healthcare regulations and institutional policies.

Fig. 2. Schema drift and robustness of learning-based semantic alignment.

### M. Summary

The proposed learning-based semantic alignment framework addresses limitations of traditional rule-based and ontology-driven methods by combining representation learning, probabilistic alignment, explainability, and governance. This architecture provides a scalable and trustworthy foundation for cross-institution health data exchange.

## V. EXPERIMENTAL METHODOLOGY

To assess the efficacy of the proposed learning-based semantic alignment framework, an experimental study was executed using realistic configurations of healthcare schemata representing cross-institutional data exchange scenarios. The experimental design evaluates robustness to heterogeneity, drift and contextually varying data, while ensuring compliance with the data governance requirements of healthcare organizations.

### A. Evaluation Objectives

The experimental methodology was designed to respond to the following research questions:

- RQ1: Is learning-based semantic alignment better than the existing rule-based and ontology-driven methods for improving mapping accuracy?
- RQ2: How robust is the proposed approach to schema drift and naming variation among institutions?
- RQ3: Will increased semantic alignment improve the accuracy of downstream analytics and reports?
- RQ4: Can the confidence of semantic alignment, along with explanations, be useful for governance and human-in-the-loop validation?

### B. Data Set Construction

Since the use of real patient data cannot be permitted for experimental testing, this study uses synthetic but realistic healthcare schemata created from typical healthcare business models.

*1) Source Domains:* Three of the most important healthcare source domains have been represented using healthcare schemata:

- Electronic Health Records (EHR): Patient demographics, encounters, diagnoses, procedures, and providers
- Claims and Billing Systems: Member identifiers, claim headers, line item information, procedure codes, diagnosis codes, and adjudicated outcomes

- Enrollment and Eligibility Systems: Coverage time frames, plan identifiers, subsidy indicators, eligibility determination

Each of these domains include multiple institutional schemata variations to simulate the heterogeneity found in real-world environments.

*2) Schemata Variants and Drift Simulation:* In order to measure robustness to schema drift, we introduce systematic drift into schemata through:

- Alias field names
- Reorganize structural (nested or flat)
- Change data types (numeric encoding strings)
- Add/Remove Optional Fields

These variations represent many of the changes caused by: Regulatory updates; Vendor migration; Workflow evolutions.

### C. Ground Truth Alignment

Semantic ground truths were created using a combination of:

- Reference mappings defined by experts
- Documented cross references to standards
- Schema derivation process

These reference mappings will be used as the basis for measuring the accuracy of semantic alignments.

### D. Baseline Methods

The proposed learning-based semantic alignment has been compared to three other methods that are commonly employed in healthcare data integration:

1) Rule-Based Matching: String similarity, naming conventions, deterministic transformation rules
2) Ontology-Driven Alignment: Mapping schema elements to controlled vocabulary, hierarchical ontologies
3) Hybrid Rule + Ontology Method: Combination of deterministic transformation rules, ontology lookup

These baseline methods represent typical industrial practices for data integration and provide a good basis for comparison.

### E. Evaluation Metrics

We assessed the quality of the alignments produced using our learning-based semantic alignment by using both intrinsic and extrinsic metrics.

*1) Alignment Accuracy Metrics:* Precision (Proportion of Correctly Predicted Alignments), Recall (Proportion of True Alignments Found), F1-Score (Harmonic Mean of Precision & Recall) were used to measure the accuracy of the semantically correct alignments at the schema level.

*2) Robustness Metrics:* We evaluated the robustness of the proposed approach to schema drift using:

- Stability Score: How consistent is the alignment when perturbed?
- Drift Sensitivity: How much does F1-score degrade with increasing amount of schema variation?

*3) Downstream Impact Metrics:* To assess the impact of the aligned schemas on downstream analytics tasks, such as:

- Reporting aggregated accuracy
- Cross-Institution Metric Consistency
- Reduction in Reconciliation Errors

Any improvements in these tasks demonstrate the practical value of the improved semantic alignment.

### F. Alignment Confidence and Explainability Assessment

Confidence scores for alignment estimates were evaluated in terms of:

- Correlation between confidence scores and alignment correctness
- Reducing False Positives via Thresholding Confidence Scores
- Effectiveness of Human-in-the-loop Review Triggered by Low Confidence Cases

Explainability Artifacts were assessed qualitatively by Domain Experts to determine if alignment decisions were actionable and interpretable.

### G. Experimental Protocol

Our evaluation follows a standard protocol:

1) Generate schema variant, generate drift scenario
2) Apply each alignment method separately
3) Produce alignment outputs and confidence scores
4) Compare output predictions with ground truth
5) Measure Performance of Downstream Analytics
6) Evaluate Governance and Explainability Artifacts

By doing so, we ensure our evaluation is reproducible and all methods are fairly compared.

### H. Summary

Overall, our experimental methodology represents a comprehensive framework for assessing the semantic alignment accuracy, robustness and governance suitability. By generating controlled schema variations within realistic healthcare environments, our evaluations can accurately model some of the real world challenges that exist in cross-institutional health data exchanges.

## VI. EXPERIMENTAL RESULTS

In this section we present experimental results of our proposed learning-based framework for semantic alignment against baseline approaches. The results are presented in terms of alignment accuracy, robustness to schema drift, downstream analytics impact, and governance-related measures such as confidence estimation and explainability.

### A. Alignment Accuracy

Alignment accuracy was measured by means of precision, recall and F1-score for each variant of schema discussed in Section VI. In Table I we summarize the average results obtained for the three domains (EHR, Claims and Enrollment).

On average our proposed solution improved the F1-score of the two traditional baselines by 20–28%. Significant gains

TABLE I
SEMANTIC ALIGNMENT ACCURACY COMPARISON

| Method | Precision | Recall | F1-Score |
| --- | --- | --- | --- |
| Rule-Based Matching | 0.71 | 0.64 | 0.67 |
| Ontology-Driven Alignment | 0.76 | 0.69 | 0.72 |
| Hybrid Rule + Ontology | 0.79 | 0.73 | 0.76 |
| Learning-Based (Proposed) | 0.88 | 0.84 | 0.86 |

were observed in those domains where there is a high degree of variability in the schema, i.e., claims and enrollment systems.

### B. Robustness to Schema Drift

We evaluated the robustness to schema drift by introducing an increasing degree of structural and naming variations in the schema. Figure 2 (already illustrated in Section V) shows some examples of drift scenarios. As shown in the figure, the degradation of rule-based methods is very rapid when the degree of drift increases; ontology-driven methods show partial resilience but struggle with structural drift. On the other hand, the learning-based framework remains stable under drift because it relies on contextual and statistical features. When the degree of drift reaches severe levels, our proposed framework retains about 82% of the baseline F1-score, whereas rule-based methods retain about 60%.

### C. Error Analysis

Error analysis has identified several common failure modes among baseline methods:

- Mismappings due to ambiguous field names
- Failures in aligning semantically equivalent but structurally distant concepts
- Over-alignments caused by overly aggressive string similarity thresholds

Our learning-based solution significantly reduces these errors by leveraging multi-modal feature representations. The remaining errors mostly involve highly specialized, institution-specific fields that provide limited contextual information.

### D. Downstream Analytics and Reporting Impact

To assess the practical relevance of our proposed solution, we performed cross-institution analytics tasks (i.e., metric aggregation and reporting consistency checks) on the aligned schemas. The key findings are:

- Reduction of 15–22% in reconciliation errors
- Consistency improvement in cross-source reports
- Decrease in manual interventions during data integration

These improvements show that the quality of semantic alignment has a direct effect on the reliability and decision-making capabilities of downstream systems.

### E. Alignment Confidence Evaluation

Alignment confidence values generated by our framework have a strong correlation with the correctness of alignments. Alignments with confidence values greater than 0.85 had a precision above 95%, while low-confidence alignments accounted

for most of the false positives. The routing of low-confidence cases to human reviewers resulted in a reduction of incorrect automatic mappings of about 40%. This demonstrates the effectiveness of the confidence-aware governance supported by our framework.

### F. Explainability Assessment

Domain experts who are responsible for data governance and interoperability have assessed the explainability artifacts provided by our framework. These experts reported that explanations based on the contributions of features (textual, statistical, contextual) are:

- Intuitive and easily understandable
- Helpful for verifying the correctness of alignment decisions
- Useful for documenting audits and compliance

Feedback from experts suggests that explainability is crucial for operational use in the context of regulated healthcare environments.

### G. Performance and Scalability Considerations

Benchmarking the performance of our framework indicates that it scales linearly with the number of schema elements. We successfully batch-aligned large schemas (over 10,000 elements) within operationally acceptable timeframes, which supports both onboarding and periodic re-alignment workflows.

### H. Summary of Findings

In summary, the experimental results presented in this paper demonstrate that our proposed learning-based framework for semantic alignment:

- Significantly improves the accuracy of alignment
- Shows a good level of resistance to schema drift and variability
- Improves the reliability of downstream analytics
- Supports governance through confidence estimation and explainability

Therefore, our findings support the applicability of our framework in real-world cross-institutional healthcare data exchange.

## VII. Conclusion

In this paper, we present a learning based framework for semantic alignment that is intended to help resolve many long-standing problems with cross-institutional healthcare data exchange. The framework combines schema representation learning, context aware similarity modeling, confidence estimation in alignments, explainability and governance to provide a scalable and trustworthy method for achieving semantic interoperability. Our experimental results demonstrate that our proposed approach substantially outperformed both traditional rule-based and ontology driven approaches in terms of alignment accuracy, robustness to schema drift and reliability of downstream analytics. In addition, the inclusion of confidence aware governance and explanations contributes to compliance

with regulations and operational trust. As healthcare systems continue to grow and connect, it will be necessary to have an adaptive and transparent method for achieving semantic alignment in order to realize the full benefit of interoperable data driven healthcare.

## References

[1] E. Tabassi, "Artificial intelligence risk management framework (ai rmf 1.0)," 2023. [Online]. Available: https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf

[2] U.S. Department of Health & Human Services, "Hipaa privacy rule," 2022. [Online]. Available: https://www.hhs.gov/hipaa/index.html

[3] Centers for Medicare & Medicaid Services, "Cms program integrity manual," 2021. [Online]. Available: https://www.cms.gov/Regulations-and-Guidance/Guidance/Manuals/Internet-Only-Manuals-IOMs-Items/CMS019033

[4] Health Level Seven International, "Hl7 version 2 product suite," 2024. [Online]. Available: https://www.hl7.org/implement/standards/product_brief.cfm?product_id=185

[5] ——, "Fhir r4 (v4.0.1): Fast healthcare interoperability resources," 2019. [Online]. Available: https://hl7.org/fhir/R4/

[6] ASC X12, "X12 standards," 2024. [Online]. Available: https://x12.org/products

[7] Health Level Seven International, "Clinical document architecture (cda)," 2024. [Online]. Available: https://www.hl7.org/implement/standards/product_brief.cfm?product_id=7

[8] SNOMED International, "Snomed ct," 2024. [Online]. Available: https://www.snomed.org/snomed-ct

[9] World Health Organization, "International classification of diseases (icd)," 2019. [Online]. Available: https://www.who.int/standards/classifications/classification-of-diseases

[10] Regenstrief Institute, "Loinc," 2024. [Online]. Available: https://loinc.org/

[11] U.S. National Library of Medicine, "Rxnorm," 2024. [Online]. Available: https://www.nlm.nih.gov/research/umls/rxnorm/

[12] E. Rahm and P. A. Bernstein, "A survey of approaches to automatic schema matching," *The VLDB Journal*, vol. 10, no. 4, pp. 334–350, 2001.

[13] P. Shvaiko and J. Euzenat, "A survey of schema-based matching approaches," in *Journal on Data Semantics IV*. Springer, 2005, pp. 146–171.

[14] A. Doan, A. Halevy, and Z. Ives, *Principles of Data Integration*. Morgan Kaufmann, 2012.

[15] D. Rodrigues, P. Ferreira *et al.*, "A study on machine learning techniques for the schema matching task," *Journal of Big Data*, vol. 8, no. 1, pp. 1–28, 2021.

[16] Y. Liu *et al.*, "Magneto: Combining small and large language models for schema matching," *Proceedings of the VLDB Endowment*, 2024, preprint PDF. [Online]. Available: https://www.vldb.org/pvldb/vol18/p2681-freire.pdf

[17] S. Wang *et al.*, "Llmatch: A unified schema matching framework with small and large language models," arXiv preprint, 2025. [Online]. Available: https://arxiv.org/abs/2507.10897