# KYUNG HEE UNIVERSITY

Department of Computer Science & Engineering, KHU, South Korea

PhD Dissertation Presentation

# Semantic Sequence Contraction and Expansion for Data Interoperability

## Fahad Ahmed Satti

Fahad.satti@oslab.khu.ac.kr

**Advisors:** Prof. Sungyoung Lee, PhD
Prof. TaeChoong Chung, PhD

# PRESENTATION AGENDA

# Background

Semantic Sequence Contraction and Expansion for Data Interoperability

## Data Interoperability

The ability with which, two or more participating systems or components can reliably <u>exchange</u> data, <u>interpret</u> it, and <u>use</u> it.

*Adapted from IEEE 610.12, HL7 and Healthcare Information Management Systems Society (HIMSS)*

## Semantic Sequence Similarity

Determine if any two given entities are <u>similar</u> or dissimilar based on their respective, hidden <u>meaning</u>.

- <u>Concept Dictionaries</u>
  (Traditional, Expert Driven)
- <u>Positional Context</u>
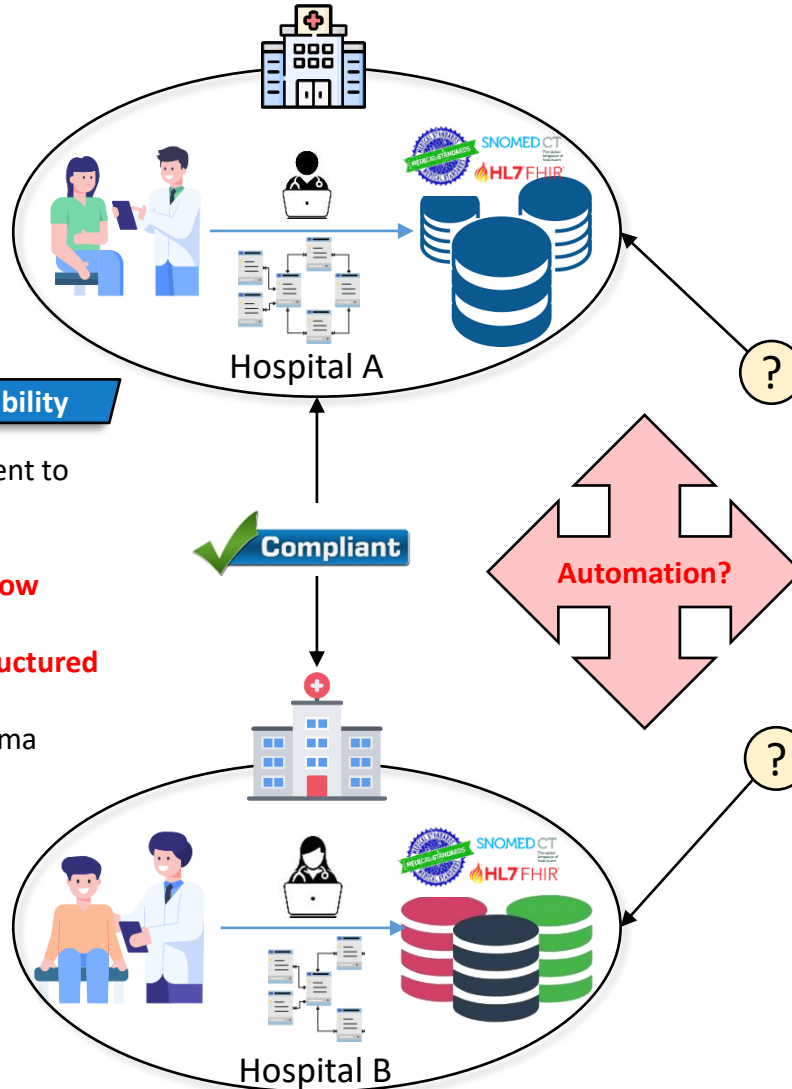  (Modern, Machine Driven)

## Standard Agnostic

A methodology, which works <u>independently</u> of any developed or under-development <u>standards</u>.

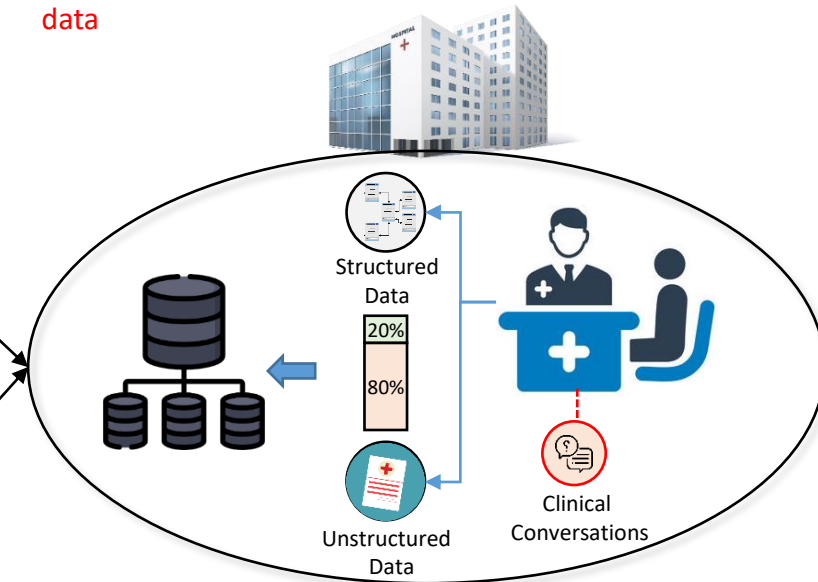*e.g. ICD-10, SNOMED-CT, LOINC, HL7 CDA, OpenEHR, HL7 FHIR*

# Motivation



**Standard-Driven Data Interoperability**

- **Expert Driven** schema alignment to exchange data between heterogeneous sources.
- Getting to an agreement is a **slow** process
- Operates on (standardized) **structured data**
- No automatic support for schema **evolution**

**Conformance for Adhoc Schema** ?

- **Resource constraints** (especially in the developing world)
- Small-Mid scale Hospitals and clinics
- Medical data with non-standard (adhoc) schema
- Unstructured data represents **80%** of medical data
- Linking unstructured data with structured medical data

## Clinical Conversations

- ❖ Primary point of data collection and inference
- ❖ Without capturing this data source
  - ❖ Some data can be lost due to cognitive load
  - ❖ Redundant effort required to digitize EMR
  - ❖ Restricts effective utilization in the developing world

# Motivation

## Standard-Driven Data Interoperability

- Expert Driven
- Highly Accurate
- Useful in the long run
- Getting to an agreement is a slow process
- Operates on structured data
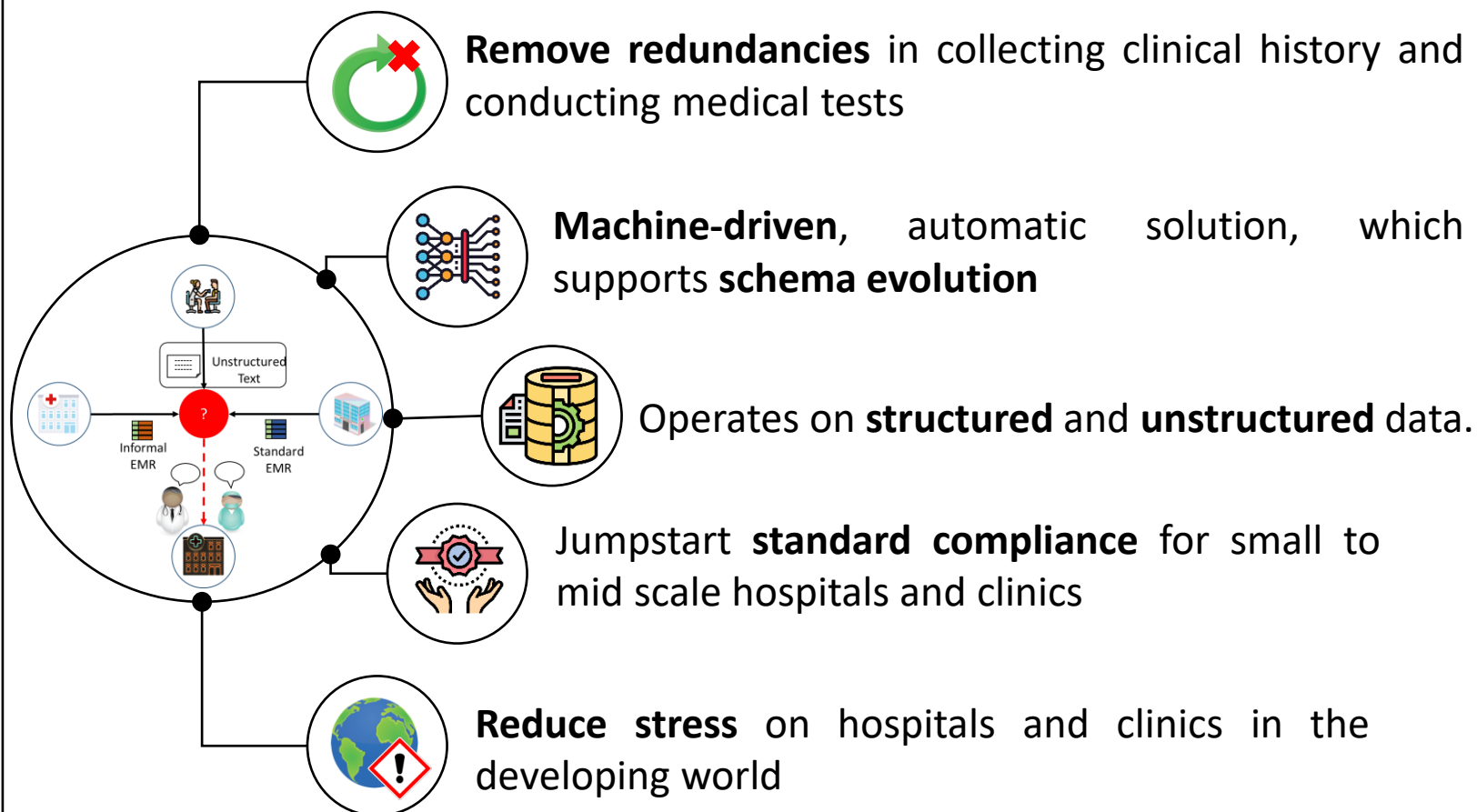- No automatic support for schema evolution

## Operational Gaps

- Difficult for small health centers to support the upgradation of existing systems
- 80% of medical data is in unstructured form (Oliver 2016)

## Active Communities

- CIMI(HL7 + OpenEHR)
- Yosemite Group
- SNOMED-CT + LOINC

## Standard-Agnostic Data interoperability

A methodology, which works independently of any formally defined standard schema or otherwise.

**Remove redundancies** in collecting clinical history and conducting medical tests

**Machine-driven**, automatic solution, which supports **schema evolution**

Operates on **structured** and **unstructured** data.

Jumpstart **standard compliance** for small to mid scale hospitals and clinics

**Reduce stress** on hospitals and clinics in the developing world

# Problem Statement

Existing Health data interoperability solutions are **expert-driven** and **standard dependent**, loosing a plethora of data residing in **informal schema** and **unstructured format**, and hindering the achievement of **Ubiquitous Healthcare**.

**Goal**

To identify and extract **clinical data** in a form consumable by various data engines for storage, usage, or exchange.

*Standard-Agnostic*

*Data Interoperability*

**Challenges**

- Challenge 1: **Identify** & **Extract** clinical **attributes** and their **values** from unstructured text
- Challenge 2: **Automatically align heterogeneous** structured and semi-structured **schema**
- Challenge 3: **Design** a **scalable** infrastructure, **automating** data interoperability.

# Research Taxonomy



[Nguyen 2019, Srihari 2008, Hara 2005, Candel 2022]

# Related Work

## Literature Survey for Sequence Contraction

| Research | Method | Advantages | Limitations |
|---|---|---|---|
| **Unstructured to Structured Data (Sequence Contraction)** Du (2019) | o Uses Bi-LSTM with CRF to first identify the sequence of interest containing a symptom and identify the target symptom<br>o Dual RNN based Seq2Seq model for identifying the similarity between utterances and existing attributes-values. | | The utterances of symptoms must be **sequential** due to the seq2seq model which relies on the in-order occurrences of symptoms. |
| Lin (2019) | o Use Bi-LSTM with a global attention mechanism to get the contextual information from document level and corpus level. The hidden layers are then re-encoded and decoded by CRF to recognize the symptoms.<br>o A symptom graph is used for symptom classification | Utilizes the semantics at document level and corpus level to identify the context of the data | o Only works on limited **pre-defined items** (authors showed results only for "upper respiratory infection", "functional dyspepsia", "infantile diarrhea" and "bronchitis").<br>o The utterances of symptoms must be **sequential** due to the use of symptom graph |
| Du (2020) | Proposed a deep learning-based approach to extract medically relevant attributes from EMR | Uses ALBERT model, which provides much better results than the traditional LSTM-CRF model. | Difficult to generalize the solution without **model retraining** |
| Zhang (2020) | Utilizes Candidate Attribute-value pairs and their status, to calculate similarity between Bert based encoded vectors for utterances and the candidates | Takes into account both the statements and question/answer type of utterances. | Only works with **existing Candidates** and is **unable to extract unseen** medical artifacts |

**Challenge 1: Limitations of existing work**
➢ Most have used **a small** set of **pre-defined attributes** which lack generalization and require intensive human efforts and time.

# Related Work

Literature Survey for Sequence Expansion

| | Research | Method | Advantages | Limitations |
|---|---|---|---|---|
| **Schema Alignment (Sequence Expansion)** | Bulygin (2018) | Devised an ontology and schema matching based approach by combining lexical and conceptual semantic similarity with various ML algorithms. | The authors have testing various ML algorithms, including Naïve Bayes, Logistic Regression, and Gradient Boosted Tree. | o Only operates on entities of **pre-defined ontologies**. <br> o All entities are matched using **naïve** comparison. |
| | Nozaki (2019) | Utilized instance-based matching and Word2Vec to create embedding vectors and calculate similarity of attributes across heterogeneous databases. | Operates on heterogeneous databases | o Word2Vec suffers from **Out of Vocabulary** problem. <br> o Only limited experiments, which do not take into account the concepts behind the values |
| | Yousfi (2020) | o Proposed an XML schema matcher, which uses conceptual semantic techniques, to transform schemas into set of words, measures each words context. <br> o Similar words are identified based on relatedness score using WordNet. | Operates on heterogeneous xml documents | o A **well-defined XSD** is necessary <br> o Only works on well formed **markup** languages <br> o **Relatedness** score of WordNet is an **old technique**, which has been replaced by the seq2seq based semantic similarity |
| | Kersloot (2020) | Reviewed several NLP algorithms for clinical text mappings onto ontological concepts. | The authors revealed that over one fourth of the NLP algorithms used were not evaluated and have no validation. | Systematic Review only |

### Challenge 2: Limitations of existing work
➢ Most solutions require a well-defined schema, which correctly and completely identifies each entity
➢ Out of vocabulary problem can greatly limit the performance of the whole technique
➢ Model trained on a specific dataset are unable to generalize

# Related Work

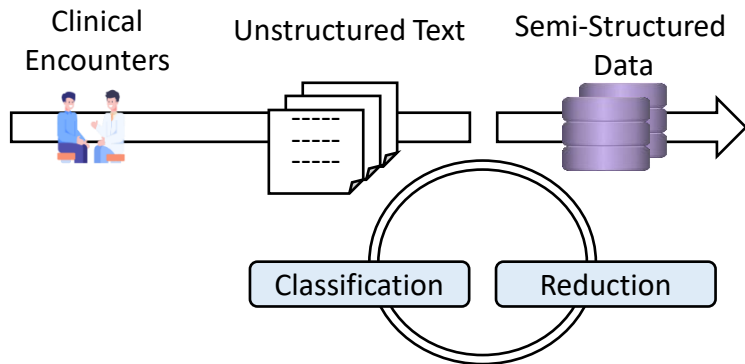## Literature Survey for Semantic Reconciliation-on-Read

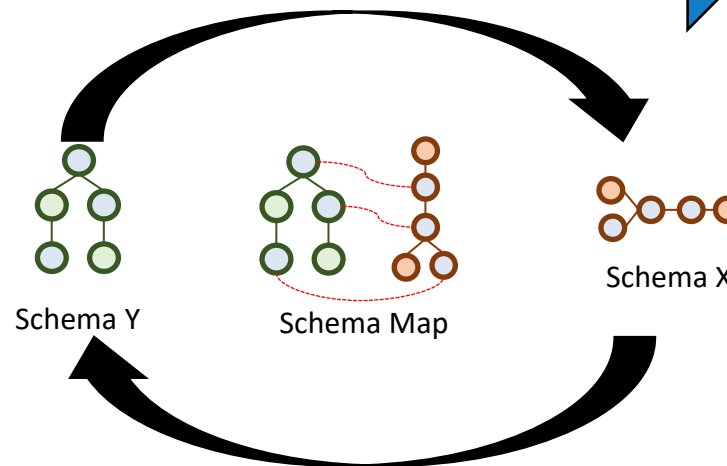| Research | Method | Advantages | Limitations |
|---|---|---|---|
| **Semantic Reconciliation-on-Read** | | | |
| LinkEHR (2019) | o Uses well defined archetypes to provide a semantic and syntactic transformation engine with large input from knowledge engineer.<br>o Depends on HL7 CDA and OpenEHR<br>o Federated query model which is based on one-to-one mapping | Provides good alignment between HL7 CDA and openEHR. | o **Standard dependent**<br>o Data retrieval dependent on how well the **transformation** definitions are.<br>o No **traceability** of healthcare records.<br>o Schema **evolution** necessitates expert input |
| OBDA (2018) | o OBDA, utilizes a well formed ontology to which all participating system must conform to.<br>o Federated query model, which does not store any data | o Does not store data, so the source data always reflects the most recent updates<br>o One to one mapping, allows any consumer or producer to provide a conformance map only once | o All systems must comply with **their standard**<br>o Data retrieval dependent on how well the **transformation** definitions are.<br>o No **traceability** of healthcare records.<br>o Schema **evolution** necessitates expert input |
| HSB (2015) | o Similar to OBDA, however the producers and consumers are loosely coupled with each other<br>o Transformation services from well-defined standard form to an internal format is required for exchanging data. | Service Bus architecture hides the details of the participating system from others | o Participating systems can comply to any system, however they should be able to **transform** the data at their ends. |

### Challenge 3: Limitations of existing work
➢ Most solutions require a well-defined schema, which correctly and completely identifies each entity
➢ No traceability of health records
➢ Schema evolution necessitates expert input to resolve any new interoperability problems
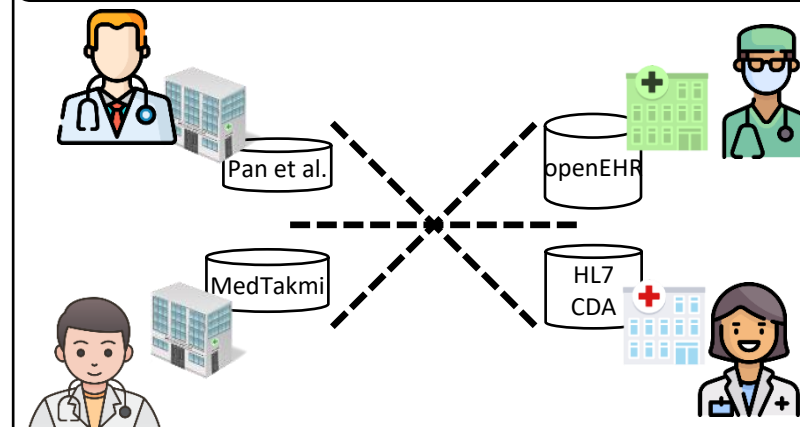
# Research Map



**Challenge 1**
**Unstructured text to semi-Structured data**

Clinical Encounters → Unstructured Text → Semi-Structured Data

Classification — Reduction

**Challenge 2**
**Alignment between schema**

Schema Y — Schema Map — Schema X

**Challenge 3**
**Semantic reconciliation-on-Read**

Pan et al.    openEHR
MedTakmi    HL7 CDA

### Existing Solutions

1. Bi-LSTM with CRF instance of interest and Dual RNN based Seq2Seq model for value identification. Du (2019), Lin (2019)
2. BERT based approaches. Zhang (2020), Du (2020)

### Limitations

o Mainly focus on a small set of attributes
o Lacks generalization.
o Require local ontologies

### Existing Solutions

1. Ontology based approach. (Bulygin 2018)
2. Instance-based matching and Word2Vec. (Nozaki 2019)
3. conceptual semantic technique working on XSD. (Yousfi 2020)

### Limitations

o Require a well-defined schema
o Out of vocabulary problem
o Lacks generalization

### Existing Solutions

1. Federated query model HL7CDA and OpenEHR. (LinkEHR 2019)
2. Federated query model, with 1-1 mapping. (OBDA 2018)
3. Health Service Bus with loose 1-1 mapping. (HSB)

### Limitations

o Require a well-defined schema
o No traceability of health records
o No support for schema evolution

# Challenges and Proposed Solutions

**Goal**

To identify and extract **clinical data** in a form consumable by various data engines for storage, usage, or exchange.

| Challenges | Proposed Solutions | Objectives |
|---|---|---|

**C1**

**Identify** & **Extract** clinical **attributes** and their **values** from unstructured text.

**S1: Sequence Contraction**

Transfer Learning to classify sequences and application of syntactic and semantic extractors for creating attribute-value pairs.

Find attributes and values from unstructured data.

**C2**

**Automatically align heterogeneous** structured and semi-structured **schema**

**S2: Sequence Expansion**

Semantic similarity of sequences, built from attribute names, using phrasal n-grams and concept enrichment.

Align Attributes with heterogeneous schema for data format transformation

**C3**

**Design** a **scalable** infrastructure, **automating** data interoperability.

**S3: Semantic Reconciliation-on-Read**

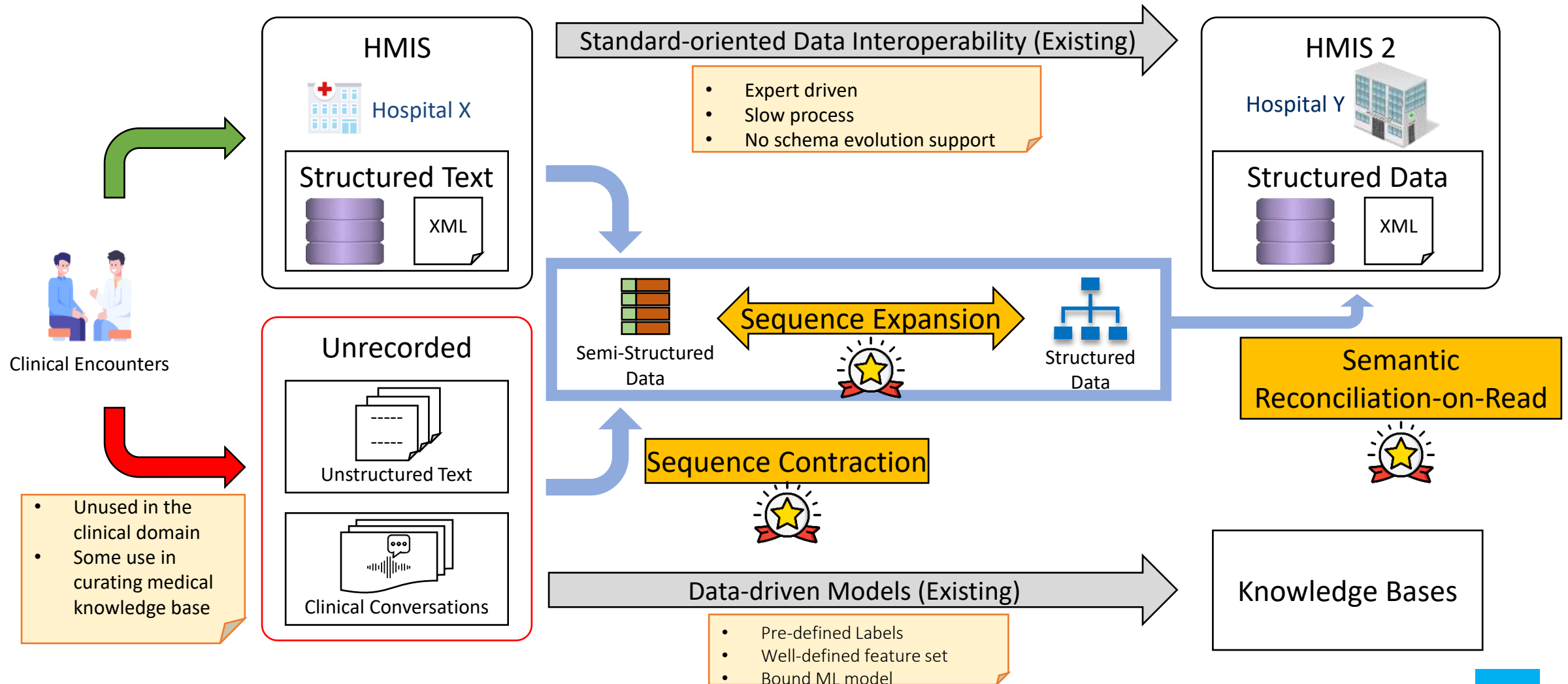A semi-structured data archiving and processing framework.

Design a practical platform which supports mapping evolution and low resource usage.

# Proposed Methodology

Idea Diagram

Key Idea
- Create Sequences from Unstructured text and attribute names
- Define a set of true sequences, enriched with semantic concepts
- Apply semantic similarity to classify unseen data
- Transform the classified instances into required results

**Clinical Encounters**

**HMIS**

Hospital X

**Structured Text**

XML

**Standard-oriented Data Interoperability (Existing)**
- Expert driven
- Slow process
- No schema evolution support

**HMIS 2**

Hospital Y

**Structured Data**

XML

**Unrecorded**

Unstructured Text

Clinical Conversations

- Unused in the clinical domain
- Some use in curating medical knowledge base

Semi-Structured Data

**Sequence Expansion**

Structured Data

**Sequence Contraction**

**Semantic Reconciliation-on-Read**

**Data-driven Models (Existing)**
- Pre-defined Labels
- Well-defined feature set
- Bound ML model

**Knowledge Bases**

# Proposed Methodology: Novelty

Semi-Structured Data →

Structured Data →

Semi-Structured Schema →

Schema Map →

## Solution 1: Sequence Contraction

**Identify relevant medical data from Clinical Conversations**



Unstructured text → Sequences → Classification ⊗ → Classified Sequences → Reduction ○ → Attr.-Value

Known Medically Aligned Sequences

RegEx | Concept Dict.

### Novelty

- Semantic Similarity based Classification
  - Transfer Learning to classify sequences
  - Easily extendable
  - Less training requirements
- Reduction
  - syntactic and semantic extractors for creating attribute-value pairs.
  - Captures syntactic artifacts like name, age, etc.
  - Utilizes conceptual semantics to enrich reduction

### Limitation

- Requires expert intervention to build the set of known medically aligned sequences and Regular expressions

## Solution 2: Sequence Expansion

**Convert attribute to sequences with suffixes and concepts**



Schema Attributes → Phrasal n-gram ○ → Sequences → Matching ⊗ → Schema Map

Suffix Arrays → Semantic Enrichment ○ →

### Novelty

- Phrasal n-gram
  - Identify hidden words within attribute names
  - Handles adhoc naming conventions
- Semantic Enrichment
  - Utilizes semantic concepts for enriched matching
- Matching
  - Unsupervised
  - m-m matching between enriched sequences

### Limitation

- No simple pathway for attributes with atomic names

## Solution 3: Semantic Reconciliation-on-Read

**Semi-structured data archiving and processing framework**



Semi-Struc. Data → Big Data Store → Collect Data ○ ← Patient Info

Schema Map → → Collect Schema ○ → Semantic Reconciliation ⊗ →

### Novelty

- Big Data Store
  - Low resource requirements at the end-nodes (hospitals/clinics)
  - Data archiving to prevent data loss
  - Evolvable schema-maps
- On demand transformation
  - Conversion to any standard
  - When required, utilize latest schema-map
  - Supports 1-1 and 1-m mappings
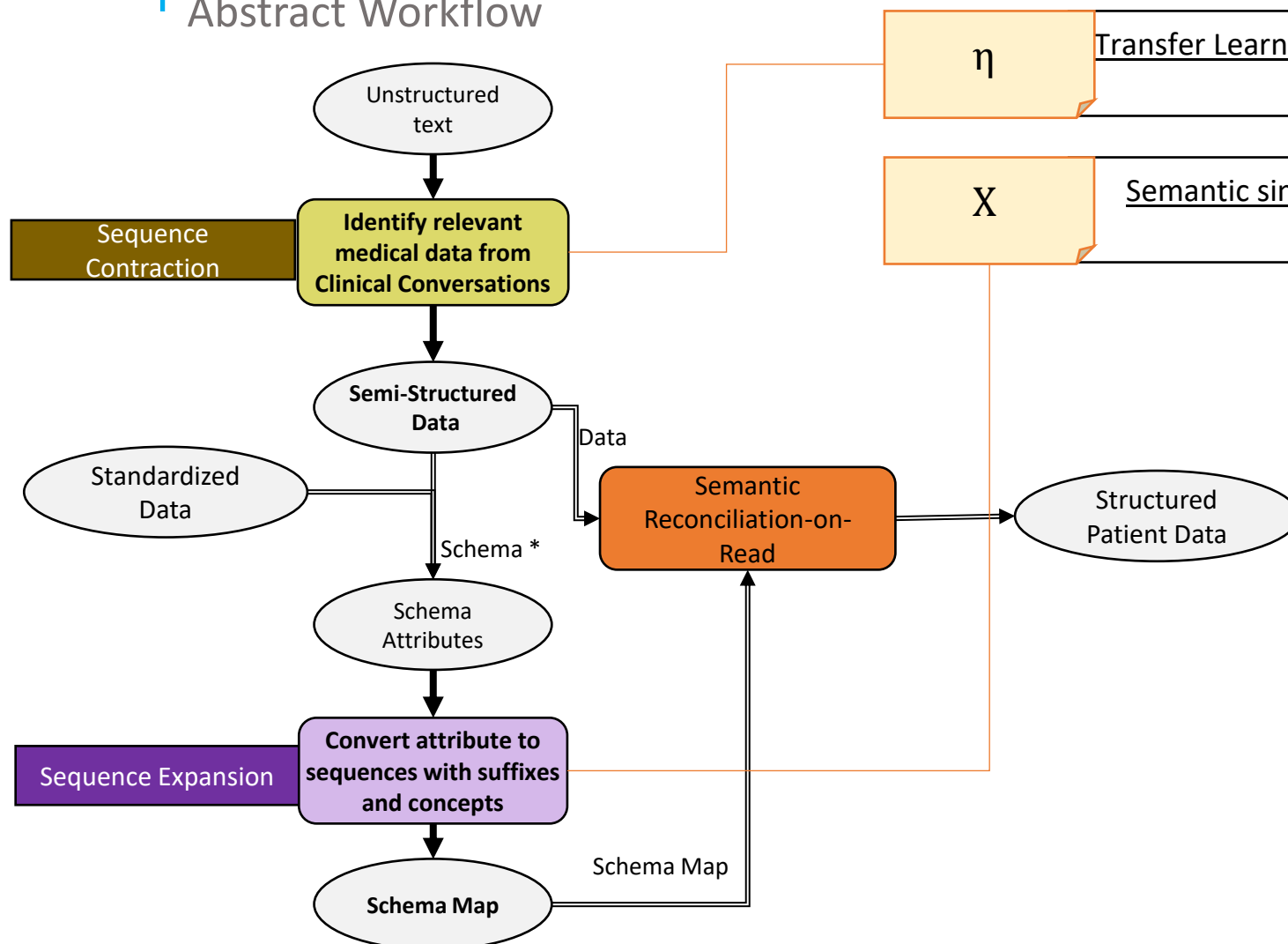
### Limitation

- Resource requirements for the Big Data Store will be high

# Proposed Methodology

## Abstract Workflow



η — Transfer Learning to classify sequences and application of syntactic and semantic extractors for creating attribute-value pairs.

X — Semantic similarity of sequences, built from attribute names, using phrasal n-grams and concept enrichment.

Unstructured text → Identify relevant medical data from Clinical Conversations

Sequence Contraction

Semi-Structured Data

Standardized Data

Data

Schema *

Schema Attributes

Semantic Reconciliation-on-Read → Structured Patient Data

Convert attribute to sequences with suffixes and concepts

Sequence Expansion

Schema Map

Schema Map

Unstructured corpus C
$$\exists C \wedge \exists \eta | \forall c \in C . \eta(c) \rightarrow p \vee \phi | p \in C$$
$$p = < p_a, p_v > | p_a \models p_v$$

Structured corpus C
$$\exists S \wedge \exists \zeta | \zeta(S) \rightarrow Q | Q = \{q\}$$
$$q = < q_a, q_v > | q_a \models q_v$$

$$\chi(p, q) = \begin{cases} 1 & if \ (p_a = q_a) \\ \sim & if \ (p_a \cong q_a) \\ 0 & otherwise \end{cases}$$

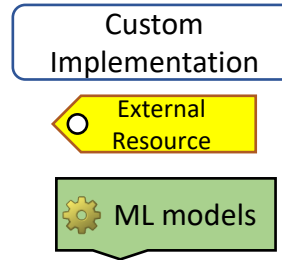* Indicates an offline and infrequent collection
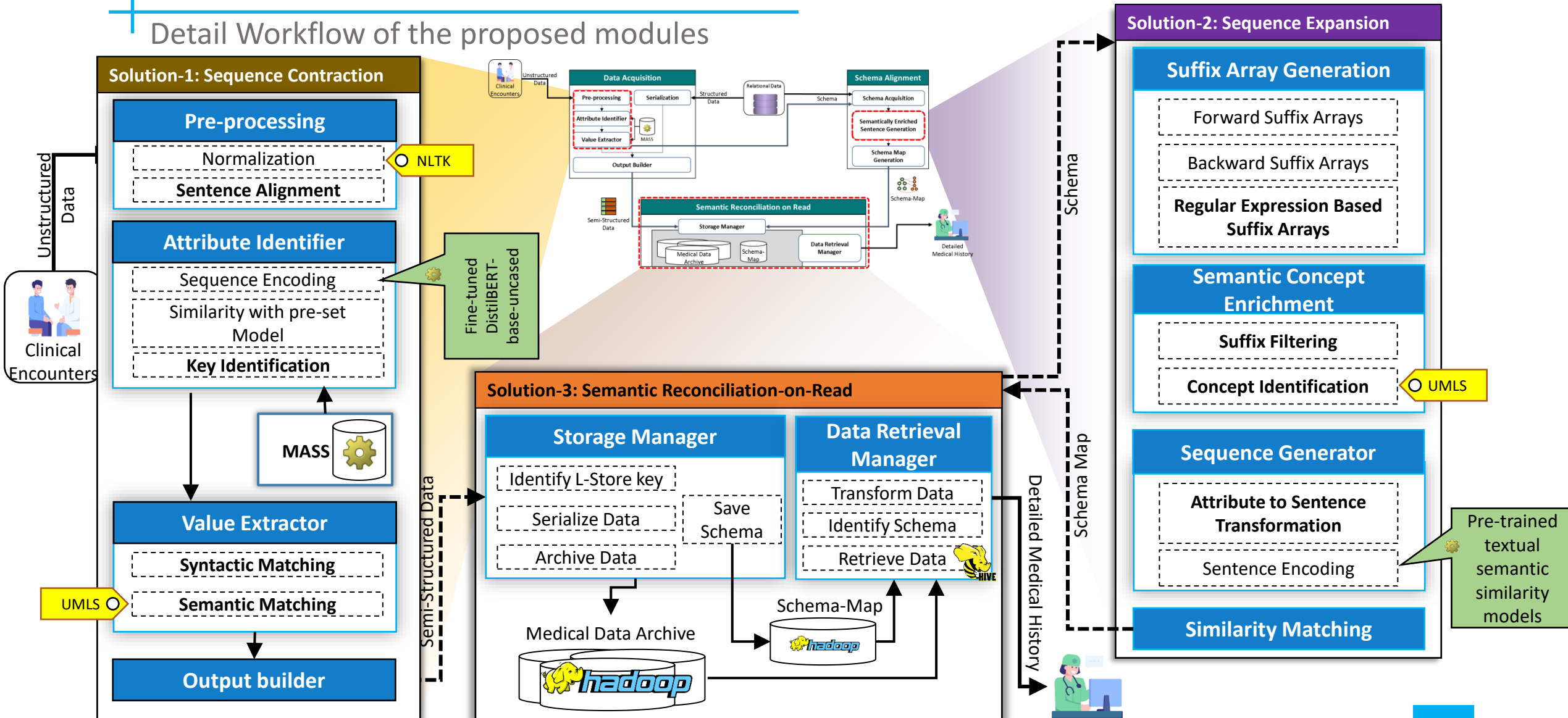
# Proposed Methodology

## Algorithmic Workflow

# Proposed Methodology

System Perspective

# Proposed Methodology

## Detail Workflow of the proposed modules



- - -→ Inter Module Communication

——→ External Communication

**Solution-1: Sequence Contraction**

### Pre-processing
- Normalization
- **Sentence Alignment**

NLTK

### Attribute Identifier
- Sequence Encoding
- Similarity with pre-set Model
- **Key Identification**

Fine-tuned DistilBERT-base-uncased

MASS

### Value Extractor
- **Syntactic Matching**
- **Semantic Matching**

UMLS

### Output builder

Clinical Encounters

Unstructured Data

**Solution-2: Sequence Expansion**

### Suffix Array Generation
- Forward Suffix Arrays
- Backward Suffix Arrays
- **Regular Expression Based Suffix Arrays**

### Semantic Concept Enrichment
- **Suffix Filtering**
- **Concept Identification**

UMLS

### Sequence Generator
- **Attribute to Sentence Transformation**
- Sentence Encoding

Pre-trained textual semantic similarity models

### Similarity Matching

Schema

Schema Map

**Solution-3: Semantic Reconciliation-on-Read**

### Storage Manager
- Identify L-Store key
- Serialize Data
- Save Schema
- Archive Data

### Data Retrieval Manager
- Transform Data
- Identify Schema
- Retrieve Data

Medical Data Archive

Schema-Map

Detailed Medical History

Semi-Structured Data

# Solution 1: Sequence Contraction

## Abstract View



Satti, Fahad Ahmed, et al. "A Semantic Sequence Similarity based approach for Extracting Medical Entities from Clinical Conversations ." *IP&M (minor revision)*

# Solution 1-1: Pre-processing

## Sequence Contraction

| AIM |
|---|
| Create sequences from text which contain both attributes and its values (Statements and Q/A) |

| Benefits |
|---|
| • Shorter self contained sequences<br>• Prioritizes shorter context over longer one<br>• Faster processing |

# Solution 1-2: Attribute Identifier

## Sequence Contraction



$$sim = \frac{\vec{V}_{S_i} \cdot \vec{V}_{\acute{S}}}{\sqrt{\vec{V}_{S_i} \cdot \vec{V}_{S_i}} \cdot \sqrt{\vec{V}_{\acute{S}} \cdot \vec{V}_{\acute{S}}}}$$

**Set of Sequences (S)**

### Model Preparation

Create the set {S x S}

Manually mark each set entry as similar or dissimilar

Expert

Create Sentence Similarity Structure ([CLS] S₁ [SEP] S₂)

Fine-tune Hyper parameters

Pre-trained DistilBERT-base-uncased Model

https://huggingface.co/distilbert-base-uncased

Fine-Tuned DistilBERT-base-uncased

### Building the Medically Aligned Sequence Set (MASS)

Expert

Annotate Sequences

| Text Sequence ($S_i$) | Generic Label ($l$) | Value Extractor ($x$) |

| Create Embedding Vector ($\vec{V_{S_i}}$) | Generic Label ($l$) | Value Extractor ($x$) |

Prepare MASS instances

MASS

### Sequence Classification

Text Sequence ($\acute{S}$)

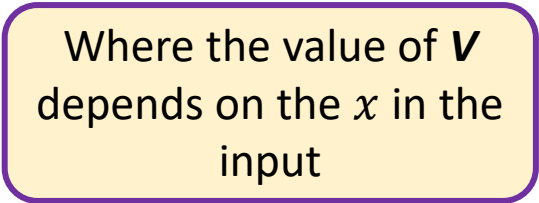Create Embedding Vector ($\vec{V_{\acute{S}}}$)

Calculate Semantic Similarity with MASS
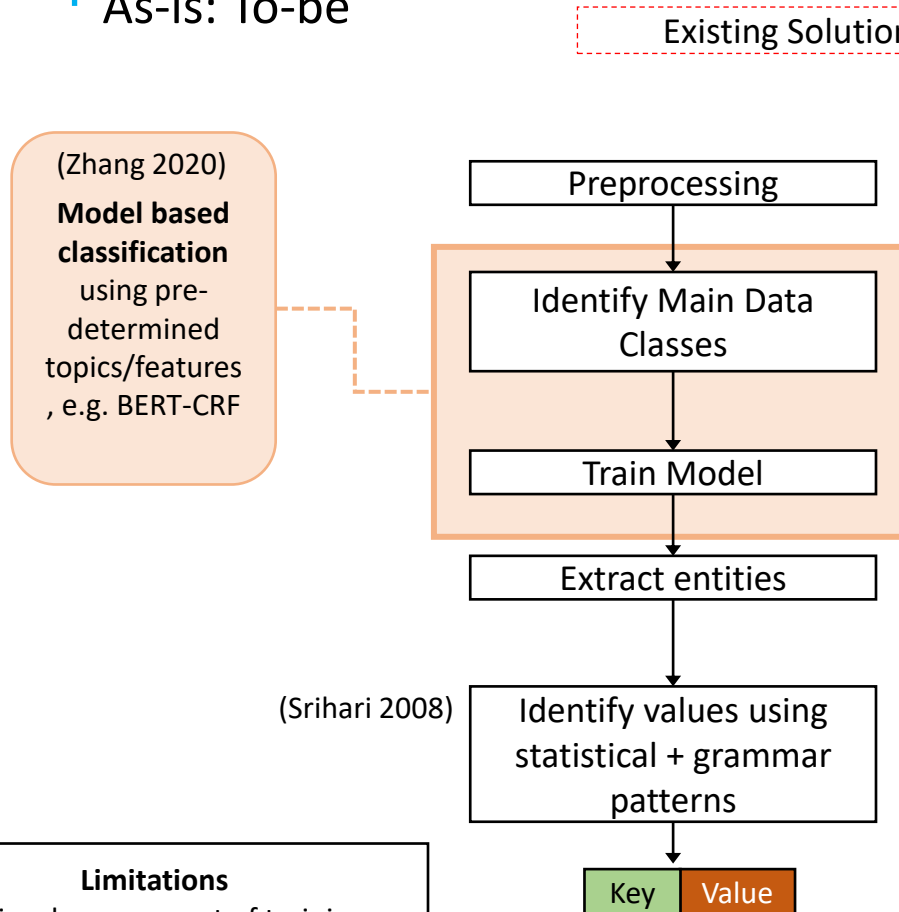
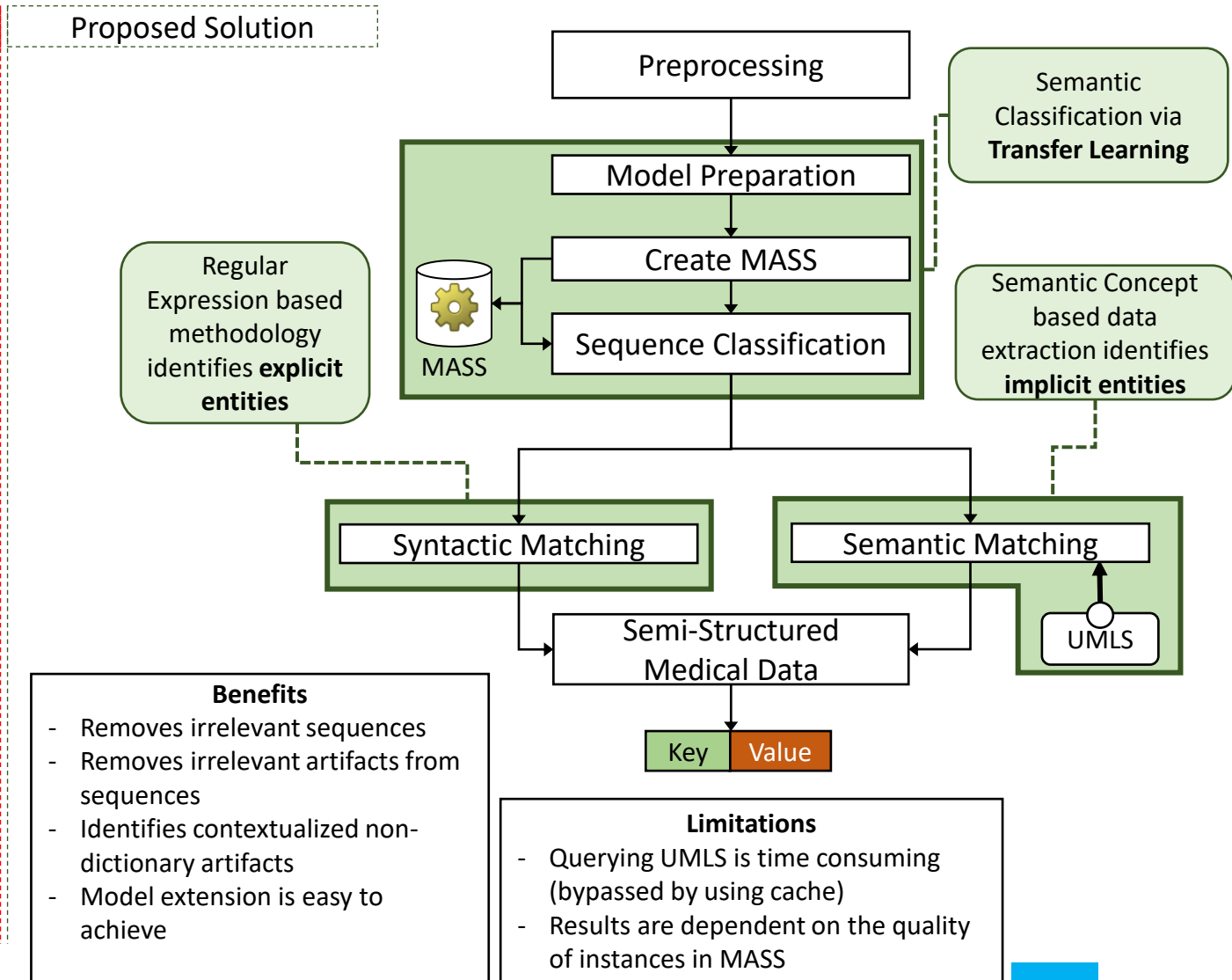Filter instances based on Threshold

| $l$ | $x$ | $\acute{S}$ |

# Solution 1-3: Value Extractor

Sequence Contraction

# Solution 1: Sequence Contraction

## As-is: To-be

Existing Solution | Proposed Solution

**(Zhang 2020)**
**Model based classification** using pre-determined topics/features , e.g. BERT-CRF

Preprocessing

Identify Main Data Classes

Train Model

Extract entities

(Srihari 2008)

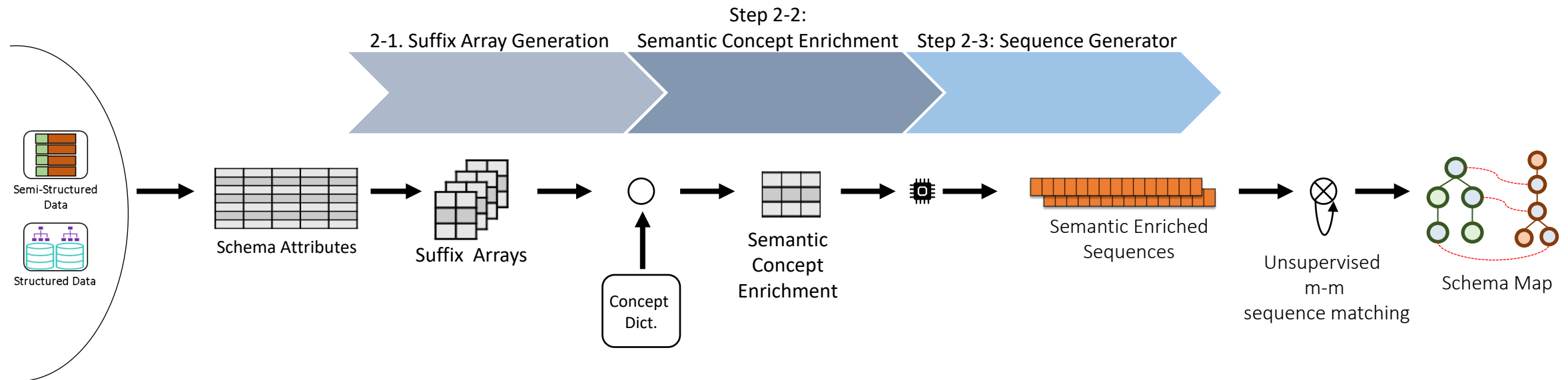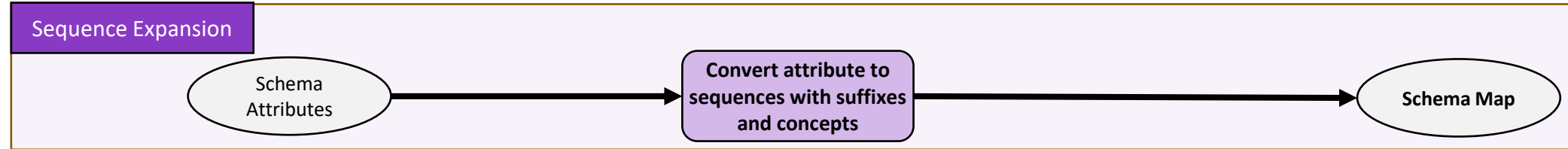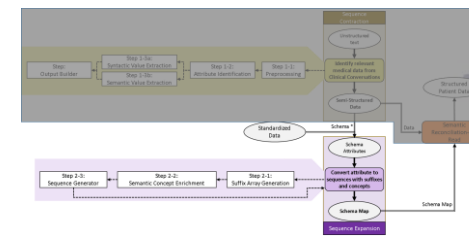Identify values using statistical + grammar patterns

Key | Value

**Limitations**
- requires large amount of training data
- Lacks generalization
- Positional semantics and feature recognition based matching only

Preprocessing

Model Preparation

Semantic Classification via **Transfer Learning**

Regular Expression based methodology identifies **explicit entities**

Create MASS

Sequence Classification

MASS

Semantic Concept based data extraction identifies **implicit entities**

Syntactic Matching

Semantic Matching

UMLS

Semi-Structured Medical Data

Key | Value

**Benefits**
- Removes irrelevant sequences
- Removes irrelevant artifacts from sequences
- Identifies contextualized non-dictionary artifacts
- Model extension is easy to achieve

**Limitations**
- Querying UMLS is time consuming (bypassed by using cache)
- Results are dependent on the quality of instances in MASS

# Solution 2: Sequence Expansion

## Abstract View



Satti, Fahad Ahmed, et al. "Unsupervised Semantic Mapping for Healthcare Data Storage Schema." *IEEE Access* 9 (2021): 107267-107278.

# Solution 2-1: Suffix Array Generation

## Sequence Expansion

| AIM | Benefits |
|---|---|
| Identify the implicit words hidden in the attribute name | • Utilizes Generalized Suffix Array; all suffixes for a set of string and is lexicographically sorted<br>• lightweight in space<br>• fast in practice |



**Algorithm 3** Suffix Array generation algorithm

1: **Input**
2:  T     token text
3: **Output**
4:  aa    Amplified Attribute
5: **procedure** BUILDSUFFIXARRAY$(T, aa)$
6:  $suffixes$: TreeSet = $\{empty\}$
7:  $N \longleftarrow length(T)$
8:  $aa : AmplifiedAttribute = \{empty\}$
9:  **for** $i \longleftarrow [1, N]$ **do**
10:  $suffixes.add(token.substring([i, N)))$
11:  **end for**
12:  **for** $j \longleftarrow [1, N)$ **do**
13:  $suffixes.add(T.substring([0, j + 1]))$
14:  **end for**
15:  $suffixes.addAll(T.split(REGEX\_WITH\_CASE))$
16:  $suffixArray$: HashSet $<$ String $> \leftarrow suffixes$
17:  **if** $suffixArray \cdot length \leq 1$ **then** return
18:  **end if**
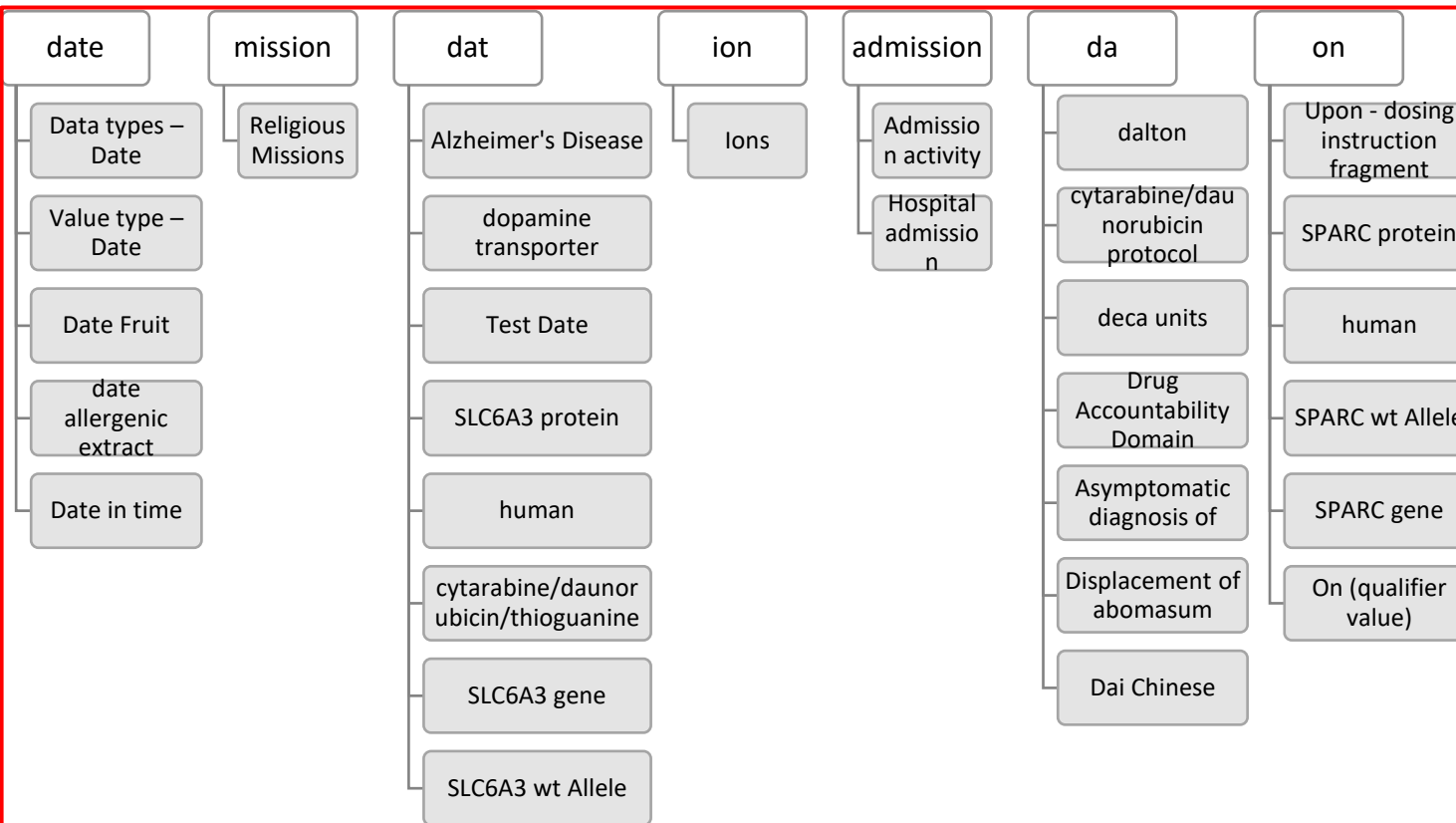19:  $aa.setSuffixes(suffixArray)$
20: **end procedure**

Forward Suffix Array

Backward Suffix Array

RegEx based Suffix Array

# Solution 2-2: Semantic Concept Enrichment

Sequence Expansion

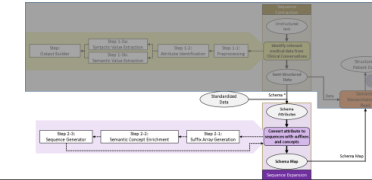| AIM | Benefits |
|---|---|
| For each suffix, identify the associated concepts | • Semantic matching can now take into account the concepts associated with each suffix |

**date**
- Data types – Date
- Value type – Date
- Date Fruit
- date allergenic extract
- Date in time

**mission**
- Religious Missions

**dat**
- Alzheimer's Disease
- dopamine transporter
- Test Date
- SLC6A3 protein
- human
- cytarabine/daunorubicin/thioguanine
- SLC6A3 gene
- SLC6A3 wt Allele

**ion**
- Ions

**admission**
- Admission activity
- Hospital admission

**da**
- dalton
- cytarabine/daunorubicin protocol
- deca units
- Drug Accountability Domain
- Asymptomatic diagnosis of
- Displacement of abomasum
- Dai Chinese

**on**
- Upon - dosing instruction fragment
- SPARC protein
- human
- SPARC wt Allele
- SPARC gene
- On (qualifier value)

**Algorithm 4** Fetch concepts from UMLS
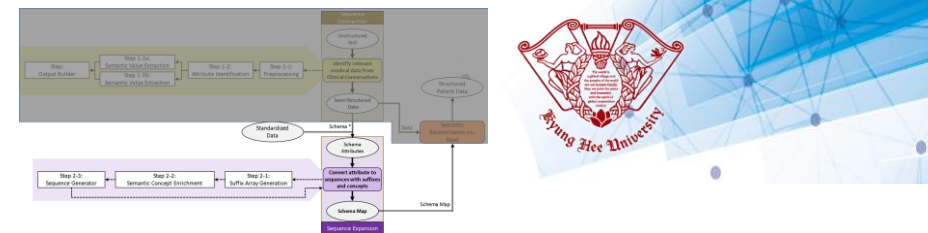
1: **Input**
2:      aa    Amplified Attribute
3: **Output**
4:      aa    Amplified Attribute
5: **procedure** FETCHUMLSCONCEPTS($aa$)
6:      $expandedTerm$: List $<$ Concept $> \leftarrow empty$
7:      **for** $word \in aa.getSuffixes$ **do**
8:         $conceptListForWord$: List $<$ Concept $> = \{empty\}$
9:         $retry \leftarrow 5$
10:         **while** $retry \neq 0$ **do**
11:            **if** $word \in umlsMap$ **then** return $umlsMap[word]$
12:            **end if**
13:            **if** $word \in umlsBlackList$ **then** return
14:            **end if**
15:            query UMLS with exact search type
16:            $results = umls.results$
17:            **for** $item \in results$ **do**
18:               **if** $item.ui = NONE$ **then** break
19:               **end if**
20:               $c_i$ : Concept $\leftarrow token, item.ui, item.name, item.root, item.uri$
21:               $expandedTerm.add(c_i)$
22:            **end for**
23:            **if** $expandedTerm = empty$ **then**
24:               $umlsBlackList.add(word)$
25:            **else**
26:               $umlsMap.put(word, expandedTerm)$
27:            **end if**
28:            **if** $exception$ **then** $retry \leftarrow retry - 1$
29:            **else** $retry \leftarrow 0$
30:            **end if**
31:         **end while**
32:         $aa.setConcepts(expandedTerm)$
33:      **end for**
34:      return aa
35: **end procedure**

# Solution 2-3: Sequence Generator

## Sequence Expansion

| AIM | Benefits |
|---|---|
| Create sequences from the amplified | • Enriches the sequence of suffixes with their concepts |

[CLS] **date** Data types - Date, Value type - Date, Date Fruit, date allergenic extract, Date in time;[SEP] **mission** Religious Missions;[SEP] **dat** Alzheimer's Disease, dopamine transporter, Test Date, SLC6A3 protein, human, cytarabine/daunorubicin/thioguanine, SLC6A3 gene, SLC6A3 wt Allele;[SEP] **ion** Ions;[SEP] **admission** Admission activity, Hospital admission;[SEP] **da** dalton, cytarabine/daunorubicin protocol, deca units, Drug Accountability Domain, Asymptomatic diagnosis of, Displacement of abomasum, Dai Chinese;[SEP] **on** Upon - dosing instruction fragment, SPARC protein, human, SPARC wt Allele, SPARC gene, On (qualifier value); [SEP]

**Algorithm 5** Create Sequences from Amplified Attribtues

1: **Input**
2:      aa    Amplified Attribute
3: **Output**
4:      $S_{exp}$ sequence expanded from a token text
5: **procedure** CREATESENTENCEFORAA($aa$)
6:      $word\_concept\_map$: Map < String, String > ← {$empty$}
7:      **for** $concept \in aa.getConcepts$ **do**
8:          $word\_concept\_map[concept.token].append(concept.name)$
9:      **end for**
10:     $S_{exp} ←$ "[CLS]"
11:     **for** $suffix \in aa.getSuffixes$ **do**
12:         $S_{exp}+ = suffix +$ " " $+$ " " $.join(word\_concept\_map[suffix]) +$ "[SEP]"
13:     **end for**
14:     $return S_{exp}$
15: **end procedure**

# Solution 2: Sequence Expansion

Detailed Workflow



$$sim = \frac{\vec{V}_{A_i} \cdot \vec{V}_{A_j}}{\sqrt{\vec{V}_{A_i} \cdot \vec{V}_{A_i}} \cdot \sqrt{\vec{V}_{A_j} \cdot \vec{V}_{A_j}}}$$

# Solution 2: Sequence Expansion

As-is: To-be



**Existing Solution**

Relational Schema

Dong (2015)
Mediated Schema

Requires expert intervention

Attribute Matching

New features in data require model retraining/fine tuning

Errica (2021)
BERT based Semantic Matching

Schema Map Generation

**Limitations**
- Based on the assumption that attributes are well defined
- Only uses the attribute for matching which produces lower accuracy

**Proposed Solution**

Relational Schema

Metadata identification

Suffix Generation

Identifies words in sequences and creates conceptual sentences

Semantic Concept Enrichment

UMLS

Sequence similarity methodology utilizing Transfer Learning

Sequence Generation

Sequence Encoding

Pre-trained Sentence Similarity Model

Sequence Matching

Schema Map Generation

**Benefits**
- Unsupervised approach
- Can deal with adhoc schemas with adhoc naming conventions
- Checks positional semantic similarity on data and its conceptual semantics, producing higher accuracy.

**Limitations**
- Unable to deal with abbreviations.

# Solution 3: Semantic Reconciliation-on-Read

## Workflow

Offline process (red arrow)
Online process (black arrow)
Remote Call (dashed arrow)

**Solution 2** → **Sequence Expansion** → Schema Map (X-Y)

**Solution 1** → **Sequence Contraction** → Semi-Structured Data (**A**)

**Semantic Reconciliation-on-Read**
- **Collect data and schema maps** as soon as they are available.
- Store them with little to **no transformation.**
- Apply the **latest schema map** on the original raw data to produce the best mappings.
- Supports mapping **evolution** and **version control** implicitly.

**Storage Manager**

Patient Identifiers (**D** + **P-Id** + **m**)     **A**

Schema Map Store

L-Store

**Medical Data Archive**

Map (X-Y)

P-Id

EHR Y — Hospital Y
Patient Medical Records

**Semi-Structured Storage Form**
- Record-Id ($i_m$)
- Type ($\tau$)
- Raw Data
- Version ($v_m$)

Medical Expert — Hospital X

Data Request (P-Id + X)

Comprehensive Medical Profile A-X/A-Y

**Data Retrieval Manager**

Integrated Data

EMR Data

| | |
|---|---|
| X-Y | Map between schema X and Y |
| A | Semi-Structured Data |
| D | Disambiguation Attributes |
| P-Id | Patient ID (UUID) |
| m | Participating Medical System |
| A-X | Data A with Schema X |

**Solution 3** → **Semantic Reconciliation on Read**

Satti, Fahad Ahmed, et al. "Ubiquitous Health Profile (UHPr): a big data curation platform for supporting health data interoperability." Computing 102.11 (2020): 2409-2444.

# Solution 3: Semantic Reconciliation-on-Read

As-is: To-be



Expert Intervention
Message Flow
Data Flow

**Existing Approaches**

Standard based Data interoperability

Federated Query based Data interoperability

**Proposed Approaches**

Semantic Reconciliation-on-Read based Data interoperability

# RESULTS AND EVALUATION

# Experimental Setup

## Solution 1: Sequence Contraction – Performance

### A sample of Test dataset

> And how old is he? 4 years;;**age:4 years**

> what happened to him?

> he has fever;;**Finding:fever**

> also some serious cough;;**Sign or Symptom:cough**

**Sequence Similarity**

1

m

Threshold Selection

### A sample of MASS instances

what is child's name? hammad

**\*[CLS] what is child's name? [MASK];;name;;(.*)?what(.*)?name(.*?)\? ((his|her|patient)? name is )?(?P<Name>.*)**

how old is he? 5 years

**\*[CLS] how old is he? [MASK] years;;age;;(old|age)(.*)?\? (he is|she is|shes)?(?P<Age>.*)(years|month)?(.*)?**

the child has cough

**\*[CLS] the child has [MASK];;Sign or Symptom;;umls**

360 Instances

79 with 1 attribute

281 with multiple attributes

|  | True Positive | True Negative |
|---|---|---|
| **Predicted Positive** | 199 | 88s |
| **Predicted Negative** | 4769 | 437 |

### Comparison



|  | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| Fine-Tuned DistilBERT | 52.96% | 69.34% | 29.44% | 41.33% |
| all-mpnet-base-v2 | 44.47% | 44.79% | 30.33% | 36.17% |
| SciBERT | 48.81% | 50.89% | 37.54% | 43.21% |

# Experimental Setup

**Solution 2:** Sequence Expansion - Dataset



| Statistics | |
|---|---|
| EMR Schema | 6 |
| Total Attributes | 270 |
| Comparisons | 48,826 |
| Annotators | 4 |

Truth Set creation by expert intervention

Cohen's Kappa score among the four annotators

A sample of 2d sheet for annotators

| | Emrbots_PatientCorePopulatedTable_PatientRace | LPanEmr_Diagnosis_heartrate |
|---|---|---|
| Openemr_Patientdata_ethnicity | 1 | 0 |
| unknown_UmlsTypes_DiagnosticProcedure | 0 | 1 |

| Annotator | Total Matches | Marked as | | | Not marked |
|---|---|---|---|---|---|
| | | Equal | Related | Unrelated | |
| Annotator 1 | 48,826 | 326 | 65+150+10 | 48275 | 0 |
| Annotator 2 | 48,826 | 329 | 36+171+25 | 48265 | 0 |
| Annotator 3 | 48,826 | 348 | 1179+884+144 | 46118 | 153 |
| Annotator 4 | 48,826 | 313 | 46+120+0 | 48336 | 11 |

# Experimental Setup

**Solution 2:** Sequence Expansion – Evaluation

Chicco (2020)

Evaluation Metric : MCC

$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP) \times (TP + FN) \times (TN + FP) \times (TN + FN)}} \rightarrow [-1,1]$$

TP = True Positive, TN = True Negative,
FP = False Positive, FN = False Negative

McHugh (2012)

Evaluation Metric : Kappa Score

$$K = \frac{(P_o - P_e)}{(1 - P_e)} \rightarrow [-1,1]$$

$P_o$ = Empirical probability of agreement on the label assigned to any sample.

$P_e$ = Expected agreement on when annotators assign labels randomly.

Annotated Dataset

Computed Matches

Agreement

**Proposed Model for Sequence Similarity**

**Existing Solutions**



| | all-mpnet-base-v2 | custom-distilbert | all-MiniLM-L6-v2 | all-distilroberta-v1 | all-MiniLM-L12-v2 | multi-qa-distilbert-cos-v1 | multi-qa-MiniLM-L6-cos-v1 | Fuzzy Matching (Senpati…) | Word2Vec (Nozaki 2019) |
|---|---|---|---|---|---|---|---|---|---|
| MCC Score | 0.43 | 0.09 | 0.32 | 0.34 | 0.06 | 0.35 | 0.33 | 0.3 | 0.02 |
| Kappa Score | 0.36 | 0.05 | 0.33 | 0.34 | 0.01 | 0.35 | 0.32 | 0.27 | 0.01 |

■ MCC Score  ■ Kappa Score

# Experimental Setup

**Solution 3:** Semantic Reconciliation-on-Read

Provides Sample data and value limits

OpenEMR [1]
12 pts.

KRSiloEMR [2]
40 pts.

EMR BOTS [3]
100k pts.

Pan et. Al [3]

MedTAKMI-CDI [4]

Data Generation →

HDFS

L-Store

**Medical Data Archive**

Semantic Query Interface

Data is loaded in 7 iterations to evaluate the scalability of the proposed data curation engine

### Final Data Statistics

| | |
|---|---|
| Total P | 390,101 |
| Total MR | 115,737,428 |

| Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 | Iteration 4 | Iteration 5 | Iteration 6 | Iteration 7 |
|---|---|---|---|---|---|---|---|
| P: 80,000 | P: 100 | P: 10,000 | P: 40,000 | P: 80,000 | P: 80,000 | P: 1 | P: 100,000 |
| MR: 2,400,000 | MR: 2,000 | MR: 200,000 | MR: 800,000 | MR: 2,400,000 | MR: 2,400,000 | MR: 40 | MR: 107,535,388 |

1. OpenEMR
2. Ali (2017)
3. Kartoun (2016)
4. Pan (2016)
5. Akihiro (2007)

P -> Generated Patients
MR -> Generated Medical Records

# Experimental Setup

**Solution 3:** Semantic Reconciliation-on-Read

| Iteration | Points | Records |
|-----------|--------|---------|
| Ite. 0 | 80,000 pts. | 2,400,000 Records |
| Ite. 1 | 100 pts. | 2,000 Records |
| Ite. 2 | 10,000 pts. | 200,000 Records |
| Ite. 3 | 40,000 pts. | 800,000 Records |
| Ite. 4 | 80,000 pts. | 2,400,000 Records |
| Ite. 5 | 80,000 pts. | 2,400,000 Records |
| Ite. 6 | 1 pt. | 40 Records |
| Ite. 7 | 100,000 pts. | 107,535,388 Records |

## Timeliness of record retreival from HDFS using Hive

Log(Time in seconds)

| | Log(C7) | Log(C8) |
|---|---------|---------|
| 1 | 1.460187965 | 2.075916867 |
| 2 | 1.454645191 | 2.084506571 |
| 3 | 1.490706957 | 2.10724729 |
| 4 | 1.518613948 | 2.143617707 |
| 5 | 1.528664788 | 2.170364115 |
| 6a | 2.096379502 | 2.288983015 |
| 6b | 1.758983008 | 2.019951659 |
| 7 | 2.423086385 | 2.991781883 |

## Scalability

Number of medical fragments / Log10 scale of medical fragments

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| C8 | 119.1 | 121.48 | 128.01 | 139.19 | 148.03 | 194.53 | 981.26 |
| Total Medical Fragments | 2E+06 | 3E+06 | 3E+06 | 6E+06 | 8E+06 | 8E+06 | 1E+08 |

# Conclusions and Future Works

## Sequence Contraction

- Proposed an automatic, semantic similarity based mechanism to extract attribute-value pairs from unstructured data

## Sequence Expansion

- Proposed a suffix array, and conceptual semantics based approach to identify the relevant parts of attribute names and used semantic similarity to align heterogeneous schemas.

## Semantic Reconciliation-on-Read

- A Big Healthare Data curation engine to archive medical data and supports schema evolution to ensure original data remains available for a longer duration

## Future Works

- The presented sequence contraction methodology can be further enhanced by increasing the sample instances in MASS.

# Publications

- **International Journals (8)**
  - First Author: 1 (Minor Revision)
  - First Author: 2 Published
  - Co-author: 5 Published
- **Local Journals (1)**
  - Co-Author: 1 Published
- **Conferences (8)**
  - First Author International: 5
  - Local Conferences: 3
- **Domestic Patents (1)**
  - Registered: 1

Publication

**Total Publications: 18**

**First Author Publications (International): 13**

# References

Oliver 2016      Müller, Oliver, et al. "Using text analytics to derive customer service management benefits from unstructured data." MIS Quarterly Executive 15.4 (2016): 243-258.

Nguyen 2019      Nguyen, D. Q., & Verspoor, K. (2019, April). End-to-end neural relation extraction using deep biaffine attention. In European conference on information retrieval (pp. 729-738). Springer, Cham.

Srihari 2008      Srihari, R. K., Li, W., Cornell, T., & Niu, C. (2008). Infoxtract: A customizable intermediate level information extraction engine. Natural Language Engineering, 14(1), 33-69.

Hara 2005      Hara, K., & Matsumoto, Y. (2005). Information Extraction and Sentence Classification applied to Clinical Trial MEDLINE Abstracts.

Candel 2022      Candel, C. J. F., Ruiz, D. S., & García-Molina, J. J. (2022). A unified metamodel for NoSQL and relational databases. Information Systems, 104, 101898.

Kartoun (2016)      Kartoun, U. (2016). A methodology to generate virtual patient repositories. arXiv preprint arXiv:1608.00570.

Pan (2016)      Pan, L., Fu, X., Cai, F., Meng, Y., & Zhang, C. (2016, October). Design a novel electronic medical record system for regional clinics and health centers in China. In 2016 2nd IEEE International Conference on Computer and Communications (ICCC) (pp. 38-41). IEEE.

Inokuchi (2007)      Inokuchi, A., Takeda, K., Inaoka, N., & Wakao, F. (2007). MedTAKMI-CDI: interactive knowledge discovery for clinical decision intelligence. IBM Systems Journal, 46(1), 115-133.

Ali (2017)      Ali, T., & Lee, S. (2017, July). Reconciliation of SNOMED CT and domain clinical model for interoperable medical knowledge creation. In 2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC) (pp. 2654-2657). IEEE.

Dong (2015)      Dong, X. L., & Srivastava, D. (2015). Schema Alignment. In Big Data Integration (pp. 31-61). Springer, Cham.

Shvaiko (2005)      Shvaiko, P., & Euzenat, J. (2005). A survey of schema-based matching approaches. In Journal on data semantics IV (pp. 146-171). Springer, Berlin, Heidelberg.

Errica (2021)      Errica, F., Silvestri, F., Edizel, B., Denoyer, L., Petroni, F., Plachouras, V., & Riedel, S. (2021, July). Concept matching for low-resource classification. In 2021 International Joint Conference on Neural Networks (IJCNN) (pp. 1-8). IEEE.

OpenEMR      OpenEMR https://www.open-emr.org/

Ali (2017)      Ali, Taqdir, and Sungyoung Lee. "Reconciliation of SNOMED CT and domain clinical model for interoperable medical knowledge creation." 2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC). IEEE, 2017.

Kartoun (2016)      Kartoun, Uri. "A methodology to generate virtual patient repositories." arXiv preprint arXiv:1608.00570 (2016).

Pan (2016)      Pan, Lijun, et al. "Design a novel electronic medical record system for regional clinics and health centers in China." 2016 2nd IEEE International Conference on Computer and Communications (ICCC). IEEE, 2016.

Akihiro (2007)      Inokuchi, Akihiro, et al. "MedTAKMI-CDI: interactive knowledge discovery for clinical decision intelligence." IBM Systems Journal 46.1 (2007): 115-133.

Chicco (2020)      Chicco, D., & Jurman, G. (2020). The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC genomics*, *21*(1), 1-13.

McHugh (2012)      Mary L McHugh. Interrater reliability: the kappa statistic. Biochemia medica, 22(3):276–282, 2012.

Zhang (2020)      Zhang, Yuanzhe, et al. "MIE: A medical information extractor towards medical dialogues." *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 2020.

# References

Lin (2019)      Lin, Xinzhu, et al. "Enhancing dialogue symptom diagnosis with global attention and symptom graph." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2019.

Du (2019)       Du, Nan, et al. "Extracting symptoms and their status from clinical conversations." *arXiv preprint arXiv:1906.02239* (2019).

Du (2020)       Du, Ming, et al. "A unified framework for attribute extraction in electronic medical records." *2020 3rd International Conference on Algorithms, Computing and Artificial Intelligence*. 2020.

Nozaki (2019)   Nozaki, Kenji, Teruhisa Hochin, and Hiroki Nomiya. "Semantic schema matching for string attribute with word vectors." *2019 6th International Conference on Computational Science/Intelligence and Applied Informatics (CSII)*. IEEE, 2019.

Yousfi (2020)   Yousfi, A., El Yazidi, M. H., & Zellou, A. (2020). xmatcher: Matching extensible markup language schemas using semantic-based techniques. *International Journal of Advanced Computer Science and Applications*, *11*(8).

Bulygin (2018)  Bulygin, L. (2018, October). Combining lexical and semantic similarity measures with machine learning approach for ontology and schema matching problem. In *Proc. Int. Conf. Data Anal. Manage. Data Intensive Domains (DAMDID/RCDL)* (pp. 245-249).

Kersloot (2020) Kersloot, M. G., van Putten, F. J., Abu-Hanna, A., Cornet, R., & Arts, D. L. (2020). Natural language processing algorithms for mapping clinical text fragments onto ontology concepts: a systematic review and recommendations for future studies. *Journal of biomedical semantics*, *11*(1), 1-21.

HSB (2015)      Meridou, Despina, et al. "An event-driven health service bus." Proceedings of the 5th EAI International Conference on Wireless Mobile Communication and Healthcare. 2015.

OBDA (2018)     Zhang, Hansi, et al. "An ontology-guided semantic data integration framework to support integrative data analysis of cancer survival." BMC medical informatics and decision making 18.2 (2018): 129-147.

Senpati (2022)  Senapati, A., Mondal, A., & Maji, S. (2022). A Fuzzy String Matching-Based Reduplication with Morphological Attributes. In *Pattern Recognition and Data Analysis with Applications* (pp. 183-190). Springer, Singapore.

# Thank you

**Comments & Suggestions**

# Prof. Seong Bae Park Comments

- **Why did you not use a classifier?**

| Classifier | Semantic Similarity |

**Training Requires a large amount of data**
- Training Requires a large amount of data
- Output labels are limited to the trained model

- Easily Extendable
- Works with small amount of pre-known data

**Classifier:**

$Label_i$

↑

$Probabilities[x_i, y_i]$

↑

Dense

↑

$Vector_i$

↑

Encoder

↑

$Token_{i1}$  $Token_{i2}$  $Token_{in}$

$Text\ Sequence_i$

**Semantic Similarity:**

Known valid Vectors   $Label_i$

↓                     ↑

Semantic Matching

↑

$Vector_i$

↑

Encoder

↑

$Token_{i1}$  $Token_{i2}$  $Token_{in}$

$Text\ Sequence_i$

https://ai.googleblog.com/2018/05/advances-in-semantic-textual-similarity.html
https://www.kaggle.com/code/thanish/bert-for-token-classification-ner-tutorial/notebook

# Prof. Eui Nam Huh Comments

- **Present an overall abstract architecture of the three solutions**

# Prof. Eui Nam Huh Comments

- **What kind of technologies are you using? (Focusing on high level architecture, how the system works, how data is collected, and how the components interact with each other)**



External Interaction
Inter Module Interaction
Intra Module Interaction

Custom Implementation
External Resource
ML models

Clinical Encounters

Unstructured Data

**Data Acquisition**

NLTK

Pre-processing

Serialization

Fine-tuned DistilBERT-base-uncased

Attribute Identifier

MASS

UMLS

Value Extractor

Output Builder

Structured Data

Relational Data

Schema

**Schema Alignment**

Schema Acquisition

Semantically Enriched Sentence Generation

UMLS

Schema Map Generation

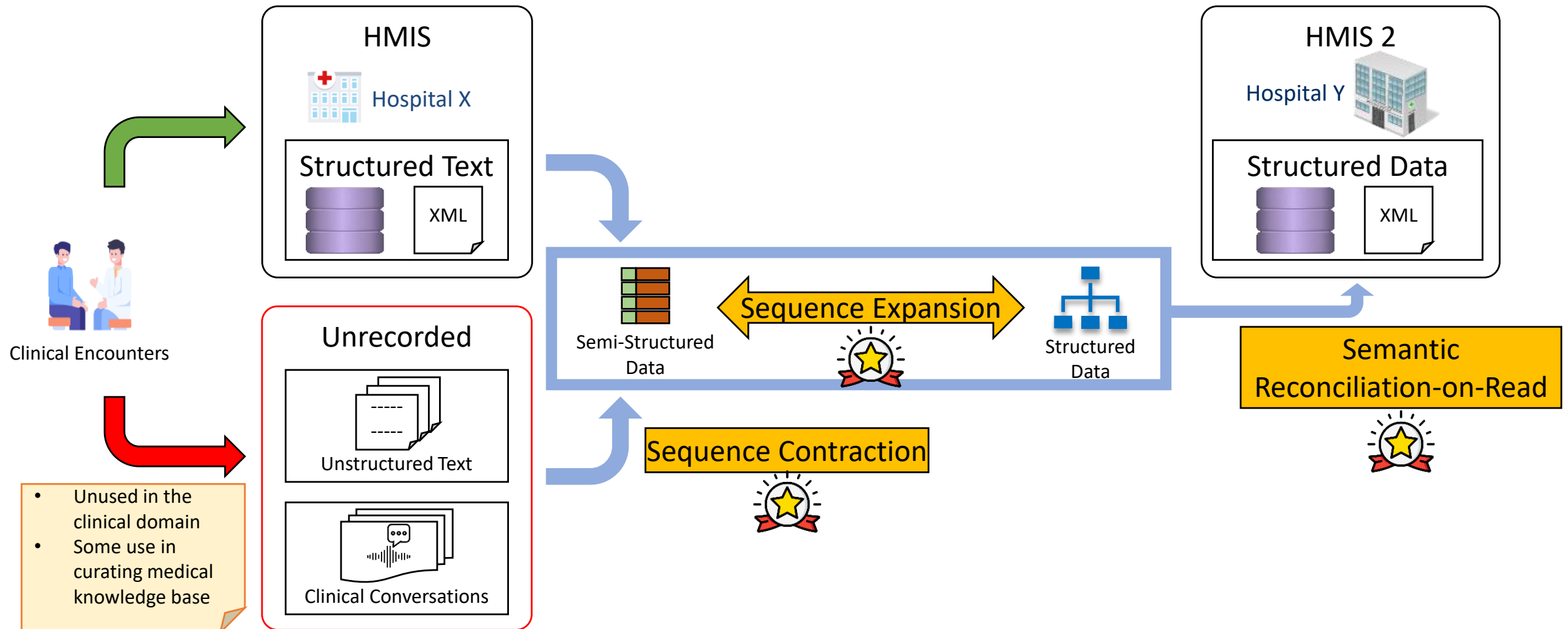Pre-trained textual semantic similarity models

S1

**Sequence Contraction**

1 x Journal Paper (IP&M)
1 x Conference Paper

S2

**Sequence Expansion**

1 x Journal Paper (IEEE Access)
3 x Conference Papers

S3

Semi-Structured Data

Schema-Map

**Semantic Reconciliation on Read**

Storage Manager

hadoop

Medical Data Archive

Schema-Map

HIVE

Data Retrieval Manager

1 x Journal Paper (Computing)
1 x Patent
2 x Conference Papers

Detailed Medical History

# Prof. Eui Nam Huh Comments

- **Why did you use Generalized Suffix Array?**

**Generalized Suffix Array:**
A generalized suffix array (or GSA), is a suffix array that contains all suffixes for a set of strings
(for example, S = S$_1$, S$_2$, S$_3$, . . . ) and is lexicographically sorted with all suffixes of each string.

The method to produce suffix array here is equivalent to the following:
S1 = dateOfAdmission
S2 = reverseElements(suffixArray(S1)) = ['on','ion','sion'.....]
S3 = ['date','Of','Admission']

S = [Admission, Of, OfAdmission, ateOfAdmission, da, dat, date, dateO,
dateOf, dateOfA, dateOfAd, dateOfAdm, dateOfAdmi, dateOfAdmis,
dateOfAdmiss, dateOfAdmissi, dateOfAdmissio, dateOfAdmission, dmission,
eOfAdmission, fAdmission, ion, ission, mission, on, sion, ssion, teOfAdmission]

**Typical GSA Usecases**: pattern matching, longest common subsequence problem, longest previous factor
array (for text compression and detection of motifs and repeats)

**Significance in Solution 2**: Use of suffix array built using only one methodology, such as Forward pass,
backward pass, or regex based one, would not capture all possible suffixes which may pertain to a medical
concept. Thus we use a combination of these three algorithms.

**Algorithm 3** Suffix Array generation algorithm
1: **Input**
2:     T     token text
3: **Output**
4:     aa     Amplified Attribute
5: **procedure** BUILDSUFFIXARRAY($T, aa$)
6:     $suffixes$: TreeSet = $\{empty\}$
7:     $N \longleftarrow length(T)$
8:     $aa : AmplifiedAttribute = \{empty\}$
9:     **for** $i \longleftarrow [1, N]$ **do**
10:        $suffixes.add(token.substring([i, N)))$
11:    **end for**
12:    **for** $j \longleftarrow [1, N]$ **do**
13:        $suffixes.add(T.substring([0, j + 1]))$
14:    **end for**
15:    $suffixes.addAll(T.split(REGEX\_WITH\_CASE))$
16:    $suffixArray$: HashSet $< String > \leftarrow suffixes$
17:    **if** $suffixArray \cdot length \leq 1$ **then** return
18:    **end if**
19:    $aa.setSuffixes(suffixArray)$
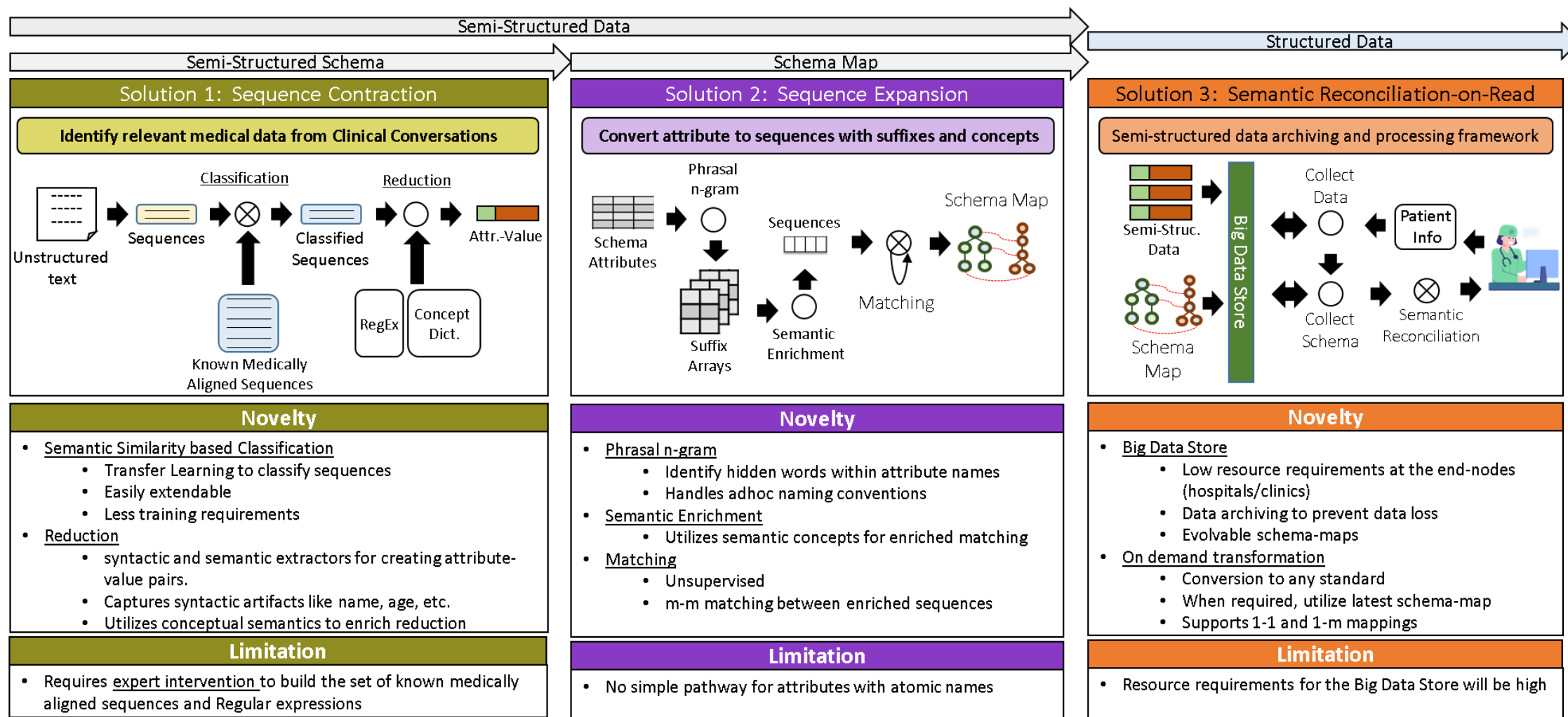20: **end procedure**

Forward Suffix Array

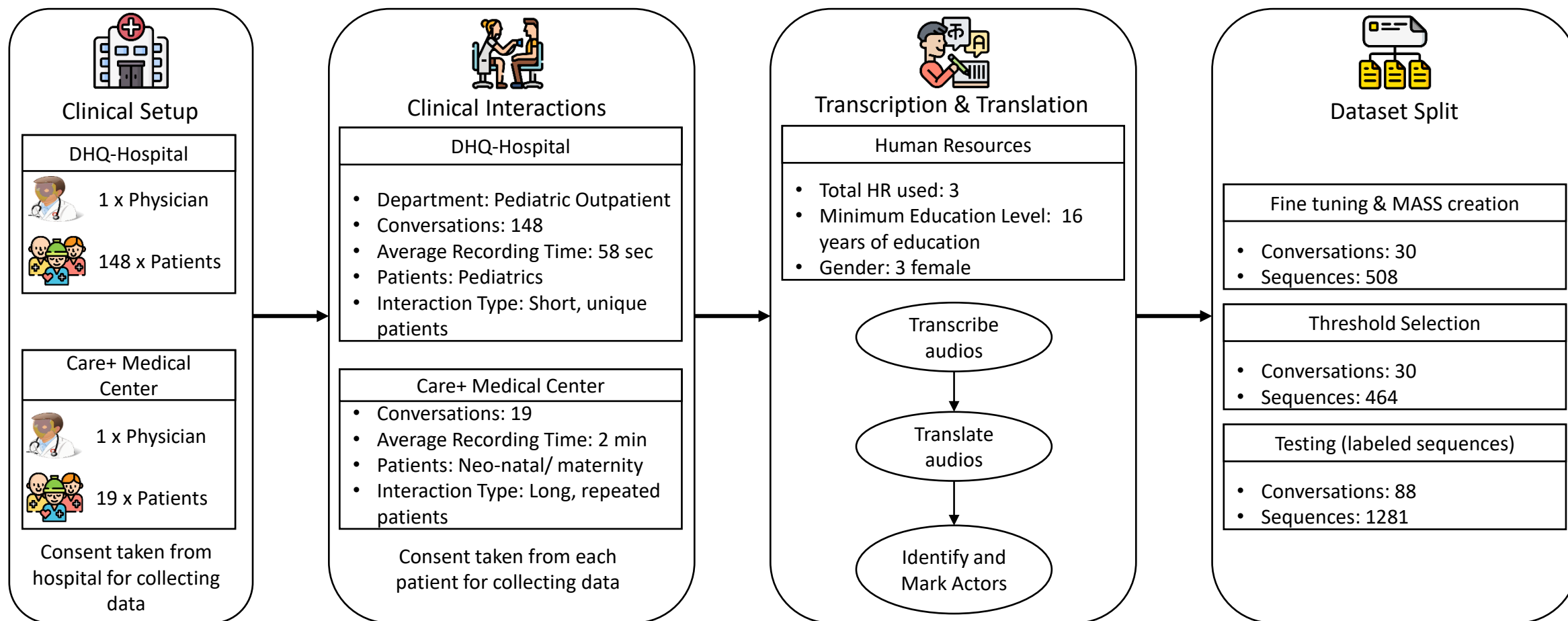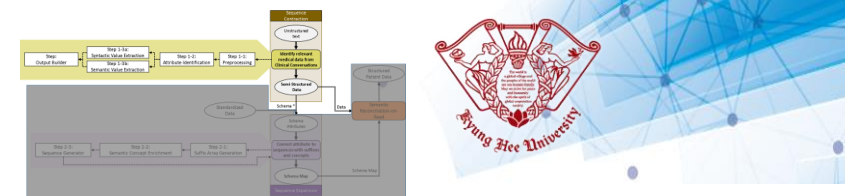Backward Suffix Array

RegEx based Suffix Array

# Prof. Eui Nam Huh Comments

- **What are your significant contributions? Focusing on strong argument for each motivation and how it is different from others.**



Semi-Structured Data → Structured Data

Semi-Structured Schema → Schema Map

## Solution 1: Sequence Contraction

**Identify relevant medical data from Clinical Conversations**

Unstructured text → Sequences → Classification → Classified Sequences → Reduction → Attr.-Value

Known Medically Aligned Sequences · RegEx · Concept Dict.

### Novelty

- Semantic Similarity based Classification
  - Transfer Learning to classify sequences
  - Easily extendable
  - Less training requirements
- Reduction
  - syntactic and semantic extractors for creating attribute-value pairs.
  - Captures syntactic artifacts like name, age, etc.
  - Utilizes conceptual semantics to enrich reduction

### Limitation

- Requires expert intervention to build the set of known medically aligned sequences and Regular expressions

## Solution 2: Sequence Expansion

**Convert attribute to sequences with suffixes and concepts**

Schema Attributes → Phrasal n-gram → Sequences → Matching → Schema Map

Suffix Arrays → Semantic Enrichment

### Novelty

- Phrasal n-gram
  - Identify hidden words within attribute names
  - Handles adhoc naming conventions
- Semantic Enrichment
  - Utilizes semantic concepts for enriched matching
- Matching
  - Unsupervised
  - m-m matching between enriched sequences

### Limitation

- No simple pathway for attributes with atomic names

## Solution 3: Semantic Reconciliation-on-Read

**Semi-structured data archiving and processing framework**

Semi-Struc. Data · Schema Map → Big Data Store → Collect Data · Collect Schema → Patient Info · Semantic Reconciliation

### Novelty

- Big Data Store
  - Low resource requirements at the end-nodes (hospitals/clinics)
  - Data archiving to prevent data loss
  - Evolvable schema-maps
- On demand transformation
  - Conversion to any standard
  - When required, utilize latest schema-map
  - Supports 1-1 and 1-m mappings

### Limitation

- Resource requirements for the Big Data Store will be high

# Experimental Setup

**Solution 1:** Sequence Contraction - Dataset

### Clinical Setup

**DHQ-Hospital**

1 x Physician

148 x Patients

**Care+ Medical Center**

1 x Physician

19 x Patients

Consent taken from hospital for collecting data

### Clinical Interactions

**DHQ-Hospital**

- Department: Pediatric Outpatient
- Conversations: 148
- Average Recording Time: 58 sec
- Patients: Pediatrics
- Interaction Type: Short, unique patients

**Care+ Medical Center**

- Conversations: 19
- Average Recording Time: 2 min
- Patients: Neo-natal/ maternity
- Interaction Type: Long, repeated patients

Consent taken from each patient for collecting data

### Transcription & Translation

**Human Resources**

- Total HR used: 3
- Minimum Education Level: 16 years of education
- Gender: 3 female

Transcribe audios

↓

Translate audios

↓

Identify and Mark Actors

### Dataset Split

**Fine tuning & MASS creation**

- Conversations: 30
- Sequences: 508

**Threshold Selection**

- Conversations: 30
- Sequences: 464

**Testing (labeled sequences)**

- Conversations: 88
- Sequences: 1281

# Solution 1: Sequence Classification

Conversation Example

## Normalization

Doctor: There's a study for which i will have to record the conversation between us regardignt the child's health, is it okay with you? Patient: Yes! Doctor: What is her name? Patient: ******. Doctor: How old is she? Patient: * months. Doctor: * months.! And what is the problem? Patient: She has temperature along with seizures. Doctor: Okay! what kind of seizures? Patient: Rapid breathing along with coughing fit. Doctor: a coughing fit along with rapid breaths and was the temerature high? Patient: Yes, its high simce yesterday. Doctor: Is she taking any feed or not? Patient: Yes she did take at 10 in the morning. Doctor: Okay, let me have a look at her, yes her respiratory rate is high. Patient: We took her to a doctor in ******, they nebulized her. Doctor: Okay, she's not fine so I'm addmitting her here, will that be okay?

Doctor: What is her name? Patient: ******.
Doctor: How old is she? Patient: * months.
Doctor: * months.! And what is the problem? Patient: She has temperature along with seizures.
Doctor: Okay! what kind of seizures? Patient: Rapid breathing along with coughing fit.
Doctor: a coughing fit along with rapid breaths and was the temerature high?
Patient: Yes, its high simce yesterday.
Doctor: Is she taking any feed or not? Patient: Yes she did take at 10 in the morning.
Doctor: Okay, let me have a look at her, yes her respiratory rate is high. Patient: We took her to a doctor in ******, they nebulized her. Doctor: Okay, she's not fine so I'm addmitting her here, will that be okay?

## Sentence Alignment

- What is her name? ******
- How old is she? * months
- And what is the problem? She has temperature along with seizures
- what kind of seizures? Rapid breathing along with coughing fit
- a coughing fit along with rapid breaths
- was the temerature high?
- its high simce yesterday
- a coughing fit along with rapid breaths and was the temerature high? Yes, its high simce yesterday
- Is she taking any feed or not? Yes she did take at 10 in the morning
- let me have a look at her
- yes her respiratory rate is high
- Okay, let me have a look at her, yes her respiratory rate is high
- We took her to a doctor in ******
- they nebulized her
- We took her to a doctor in ******, they nebulized her
- she's not fine so I'm addmitting her here
- will that be okay?
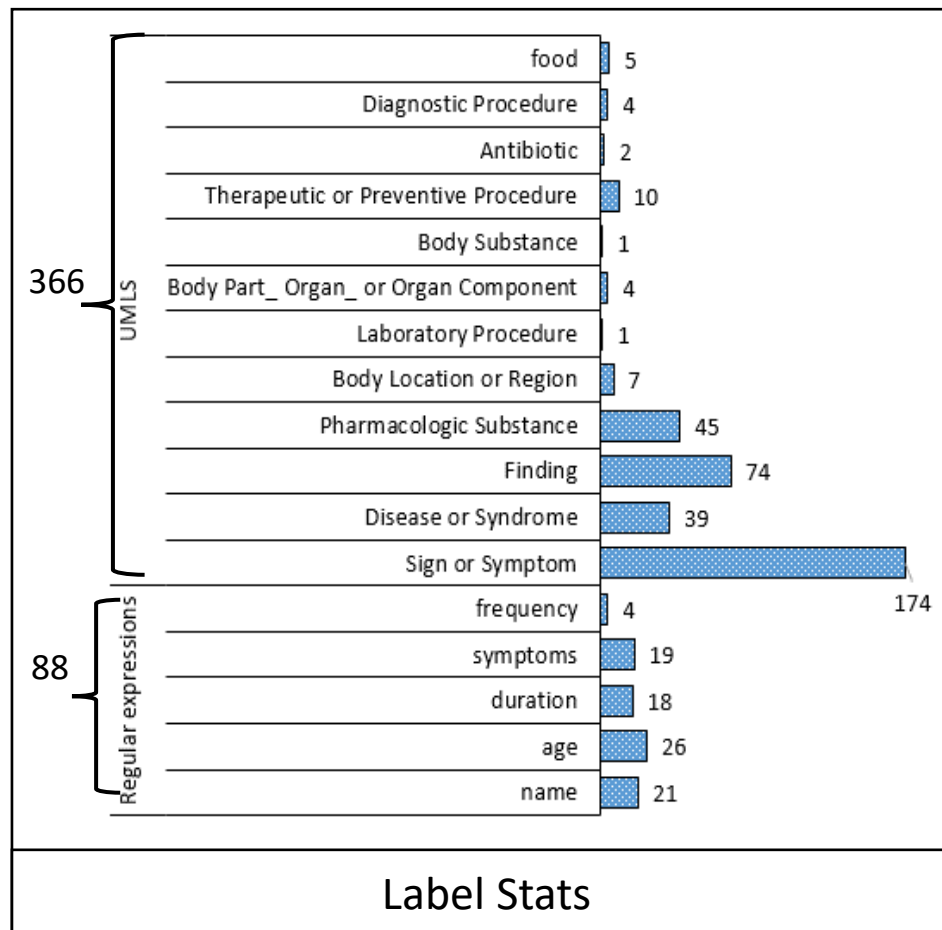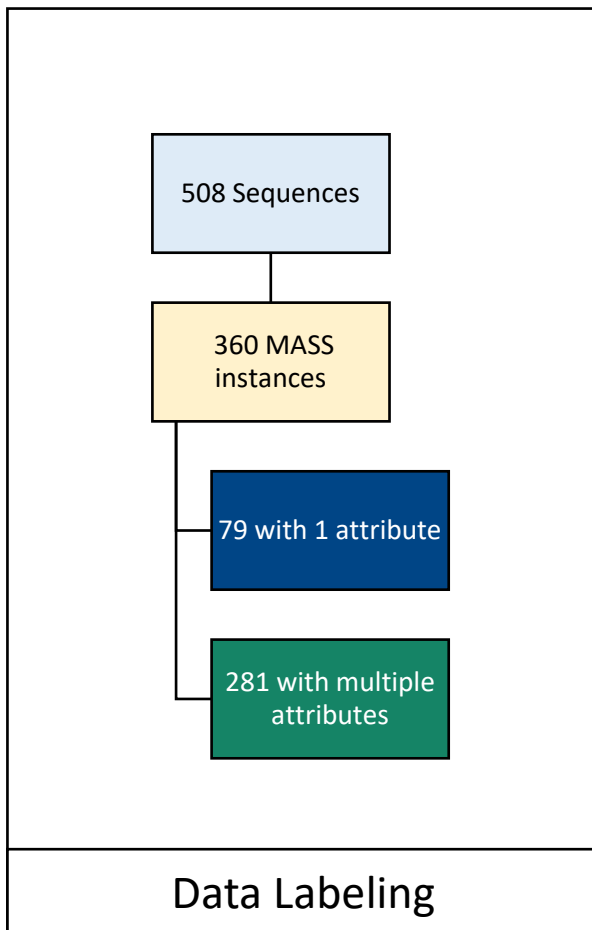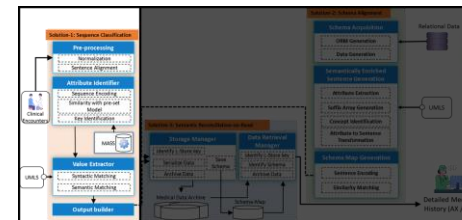- Okay, she's not fine so I'm addmitting her here, will that be okay?

# Solution 1: Fine Tuning DistilBERT



Training Loss and Accuracy on Transcripts

| Hyper parameter | Value |
|---|---|
| Batch Size | 32 |
| Loss Function | Cross Entropy |
| Evaluation Metric | Sparse Categorical Accuracy |
| Optimizer | AdamW |
| Initial Learning Rate | 1e-4 |
| Warmup steps | 10% |
| Test Accuracy | 95% |

# Experimental Setup

## Solution 1: Sequence Contraction - MASS development



**Data Labeling**

- 508 Sequences
- 360 MASS instances
- 79 with 1 attribute
- 281 with multiple attributes

**Label Stats**

UMLS (366):
- food — 5
- Diagnostic Procedure — 4
- Antibiotic — 2
- Therapeutic or Preventive Procedure — 10
- Body Substance — 1
- Body Part_ Organ_ or Organ Component — 4
- Laboratory Procedure — 1
- Body Location or Region — 7
- Pharmacologic Substance — 45
- Finding — 74
- Disease or Syndrome — 39
- Sign or Symptom — 174

Regular expressions (88):
- frequency — 4
- symptoms — 19
- duration — 18
- age — 26
- name — 21

**A sample of MASS instances**

what is child's name? hammad
*[CLS] what is child's name? [MASK];;name;;(.*)?what(.*)?name(.*?)\? ((his|her|patient)? name is )?(?P<Name>.*)

how old is he? 5 years
*[CLS] how old is he? [MASK] years;;age;;(old|age)(.*)?\? (he is|she is|shes)?(?P<Age>.*)(years|month)?(.*)?

the child has cough
*[CLS] the child has [MASK];;Sign or Symptom;;umls

what's wrong with the baby? the child has cough and cold
*[CLS]what's wrong with the baby? [SEP] the child has [MASK] and [MASK];;Sign or Symptom,Sign or Symptom;;umls

how long? it's been 3 days
*[CLS] how long? [SEP] it's been [MASK] [days];;duration;;(.* )? it's been (\s+the last|\s+previous)?\s+(?P<Duration>.*(day|week|month|year)?(s)?)

- Due to the presence of multiple attributes in MASS instances, the number of labels are more than the number of sequences
- UMLS labels correspond to the UMLS semantic concept types

# Solution 1: Sequence Classification

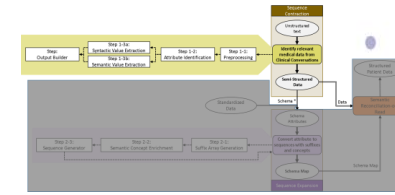Pre-processing Example

## Clinical Text

Doctor: There's a study for which i will have to record the conversation between us regardignt the child's health, is it okay with you? Patient: Yes! Doctor: What is her name? Patient: ******. Doctor: How old is she? Patient: * months. Doctor: * months.! And what is the problem? Patient: She has temperature along with seizures. Doctor: Okay! what kind of seizures? Patient: Rapid breathing along with coughing fit. Doctor: a coughing fit along with rapid breaths and was the temerature high? Patient: Yes, its high simce yesterday. Doctor: Is she taking any feed or not? Patient: Yes she did take at 10 in the morning. Doctor: Okay, let me have a look at her, yes her respiratory rate is high. Patient: We took her to a doctor in ******, they nebulized her. Doctor: Okay, she's not fine so I'm addmitting her here, will that be okay?

## Sequences

- What is her name? ******
- How old is she? * months
- And what is the problem? She has temperature along with seizures
- what kind of seizures? Rapid breathing along with coughing fit
- a coughing fit along with rapid breaths
- was the temerature high?
- its high simce yesterday
- a coughing fit along with rapid breaths and was the temerature high? Yes, its high simce yesterday
- Is she taking any feed or not? Yes she did take at 10 in the morning
- let me have a look at her
- yes her respiratory rate is high
- Okay, let me have a look at her, yes her respiratory rate is high
- We took her to a doctor in ******
- they nebulized her
- We took her to a doctor in ******, they nebulized her
- she's not fine so I'm addmitting her here
- will that be okay?
- Okay, she's not fine so I'm addmitting her here, will that be okay?
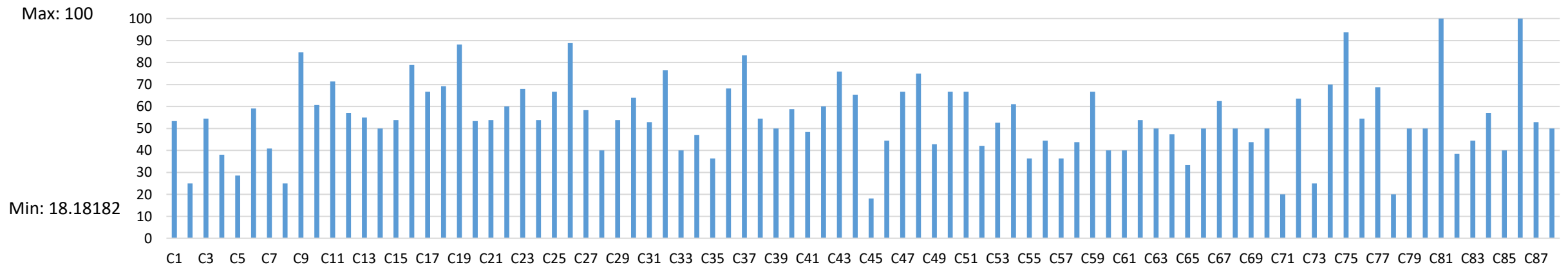
# Experimental Setup

**Solution 1:** Sequence Contraction – Individual conversation Performance
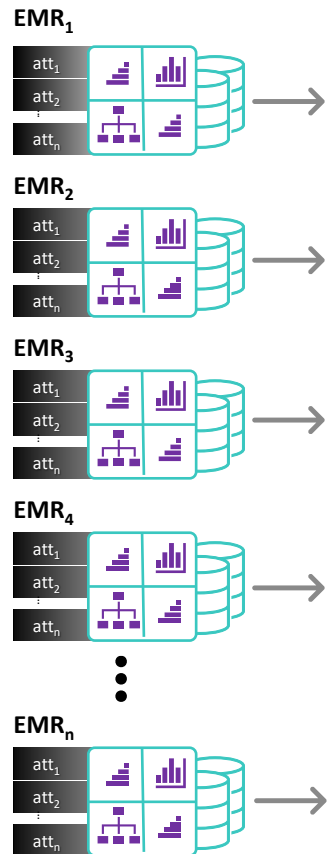
(a) **Baseline Methodology**: all-mpnet-base-v2

Max: 93.75

Min: 18.18182



(b) **Proposed Methodology**: Fine-Tuned DistilBERT

Max: 100

Min: 18.18182

# Solution 2: Sequence Expansion

## Abstract View



**Sequence Expansion**

Relational Data → Schema Acquisition → Semantically Enriched Sentence Generation (Sequence Expansion) → Schema Map Generation

EMR$_1$, EMR$_2$, EMR$_3$, EMR$_4$ ... EMR$_n$ (att$_1$, att$_2$, att$_n$)

EMR Schema → Data Instances ← Terminological Standard Dictionaries (UMLS, SNOMED-CT or LONICS)

Schema Attribute → Suffix Arrays → Semantic Concept Enrichment → Sequence Generator → Resultant Sequence

Semantic Enriched Sequences → Context vectors$_1$ (0001000000000), Context vectors$_2$ (0000001000000) → Similarity Matching → Embedded Vectors Generations → Schema-Map

To Solution 3 → Store in DB

Satti, Fahad Ahmed, et al. "Unsupervised Semantic Mapping for Healthcare Data Storage Schema." IEEE Access 9 (2021): 107267-107278.

# Solution 2: Sequence Expansion

## Challenge: Heterogeneous Schema



| | | |
|---|---|---|
| S1 | OpenEMR | https://www.open-emr.org/ |
| S2 | EMRBOTS | Kartoun (2016) |
| S3 | Lpan EMR | Pan (2016) |
| S4 | MedTAKMI-CDI EMR | Inokuchi (2007) |
| S5 | KrSilo EMR | Ali (2017) |

# Experimental Setup

**Solution 2:** Sequence Expansion – Threshold Selection

Evaluation Metric : MCC

$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP) \times (TP + FN) \times (TN + FP) \times (TN + FN)}} \rightarrow [-1,1]$$

TP = True Positive, TN = True Negative,
FP = False Positive, FN = False Negative

- Accuracy fails to account for imbalanced datasets
- F1 measure is not affected by the true negative scores.
- MCC provides an acceptable alternate in our current scenario comprising of imbalanced dataset (largely in favour of class "unrelated")
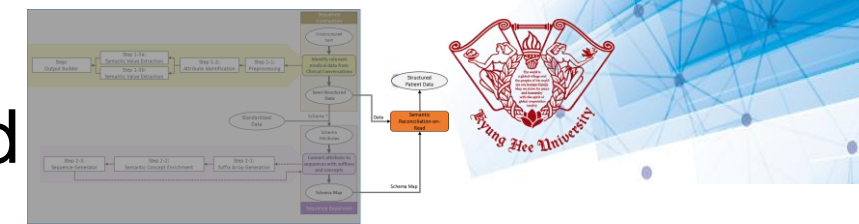
Chicco (2020)



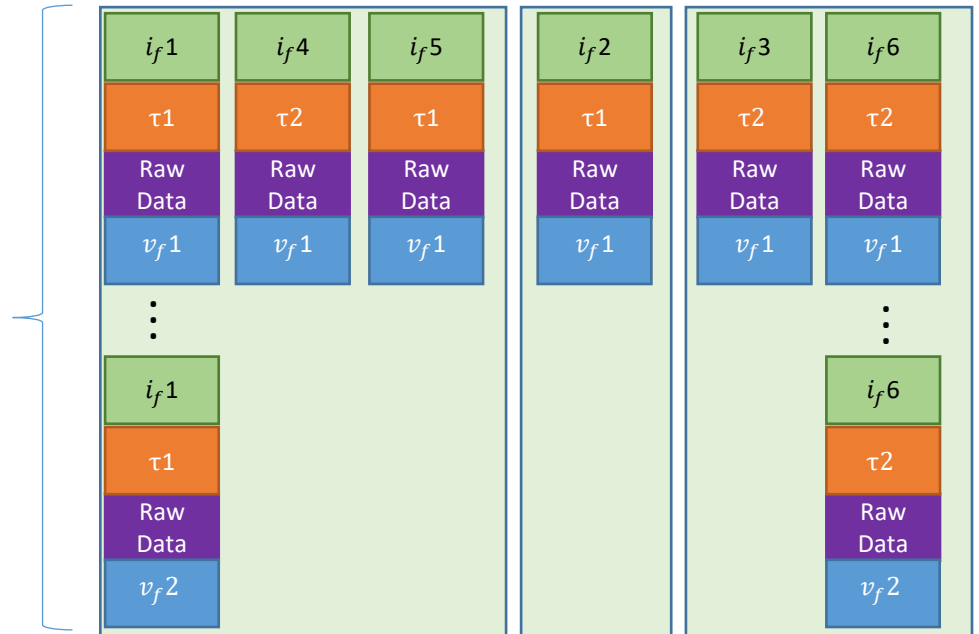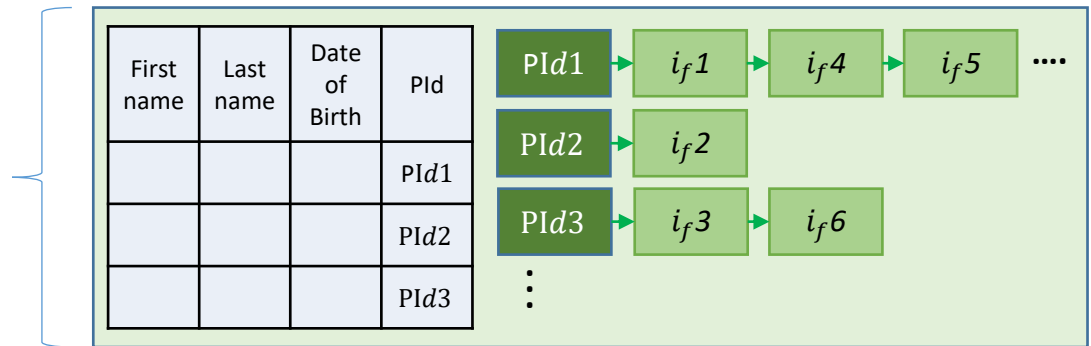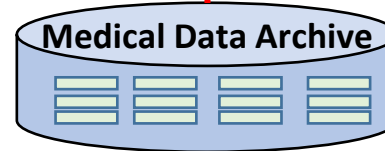| Method | Threshold | |
| --- | --- | --- |
| | Unequal and similar | Similar and equal |
| all-mpnet-base-v2 | 0.75 | 0.8 |
| custom-distilbert | 0.75 | 0.9 |
| all-MiniLM-L6-v2 | 0.55 | 0.6 |
| all-distilroberta-v1 | 0.65 | 0.8 |
| all-MiniLM-L12-v2 | 0 | 0.3 |
| multi-qa-distilbert-cos-v1 | 0.6 | 0.65 |
| multi-qa-MiniLM-L6-cos-v1 | 0.65 | 0.7 |

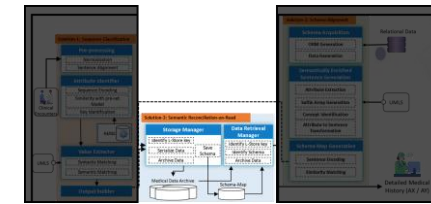# Solution 3: Semantic Reconciliation-on-Read

## Data Storage View



- EMR data is **serialized** into Semi-Structured form with "Raw Data" containing "**key:value**" pairs.
- **Disambiguation attributes** (such as firstname, lastname, dateofbirth) can be used to **identify patients** across organizational boundaries.
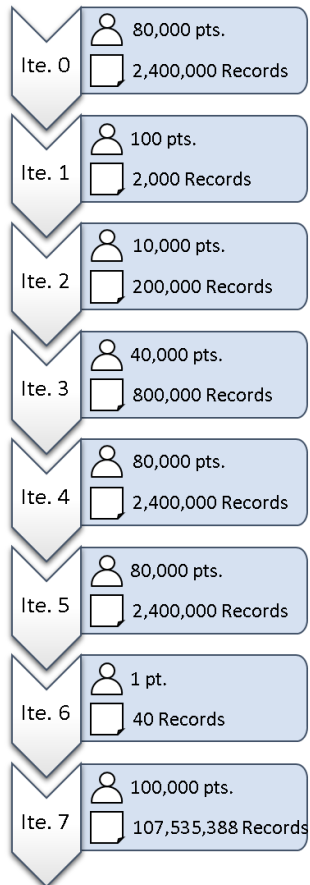
EMR Semi-structured Storage Form

# Experimental Setup

**Solution 3:** Evaluation Critera

Timeliness Evaluation for all 7 iterations shows the performance of the proposed approach in the presence of Big Data.
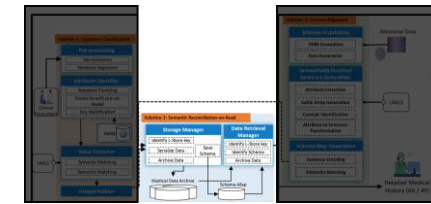
| Ite. 0 | 80,000 pts. / 2,400,000 Records |
| Ite. 1 | 100 pts. / 2,000 Records |
| Ite. 2 | 10,000 pts. / 200,000 Records |
| Ite. 3 | 40,000 pts. / 800,000 Records |
| Ite. 4 | 80,000 pts. / 2,400,000 Records |
| Ite. 5 | 80,000 pts. / 2,400,000 Records |
| Ite. 6 | 1 pt. / 40 Records |
| Ite. 7 | 100,000 pts. / 107,535,388 Records |

| Iteration | Total Fragments | File size for C1 (Kb) | File size for C2 (Kb) | File size for C3 (Kb) |
|---|---|---|---|---|
| Initial | 2,400,000 | - | - | - |
| 1 | 2000 | 659 | 6 | 181 |
| 2 | 200,000 | 66,260 | 580 | 18,059 |
| 3 | 800,000 | 264,923 | 2320 | 72,242 |
| 4 | 2,400,000 | 755,295 | 4,639 | 216,617 |
| 5 | 2,400,000 | 755,417 | 4,639 | 216,608 |
| 6 | 40 | 13 | 1 | 4 |
| 7 | 107,535,388 | 25,752,400 | 7,263 | 11,118,380 |
| **Total** | **115,737,428** | **27,594,967** | **19,448** | **11,642,091** |

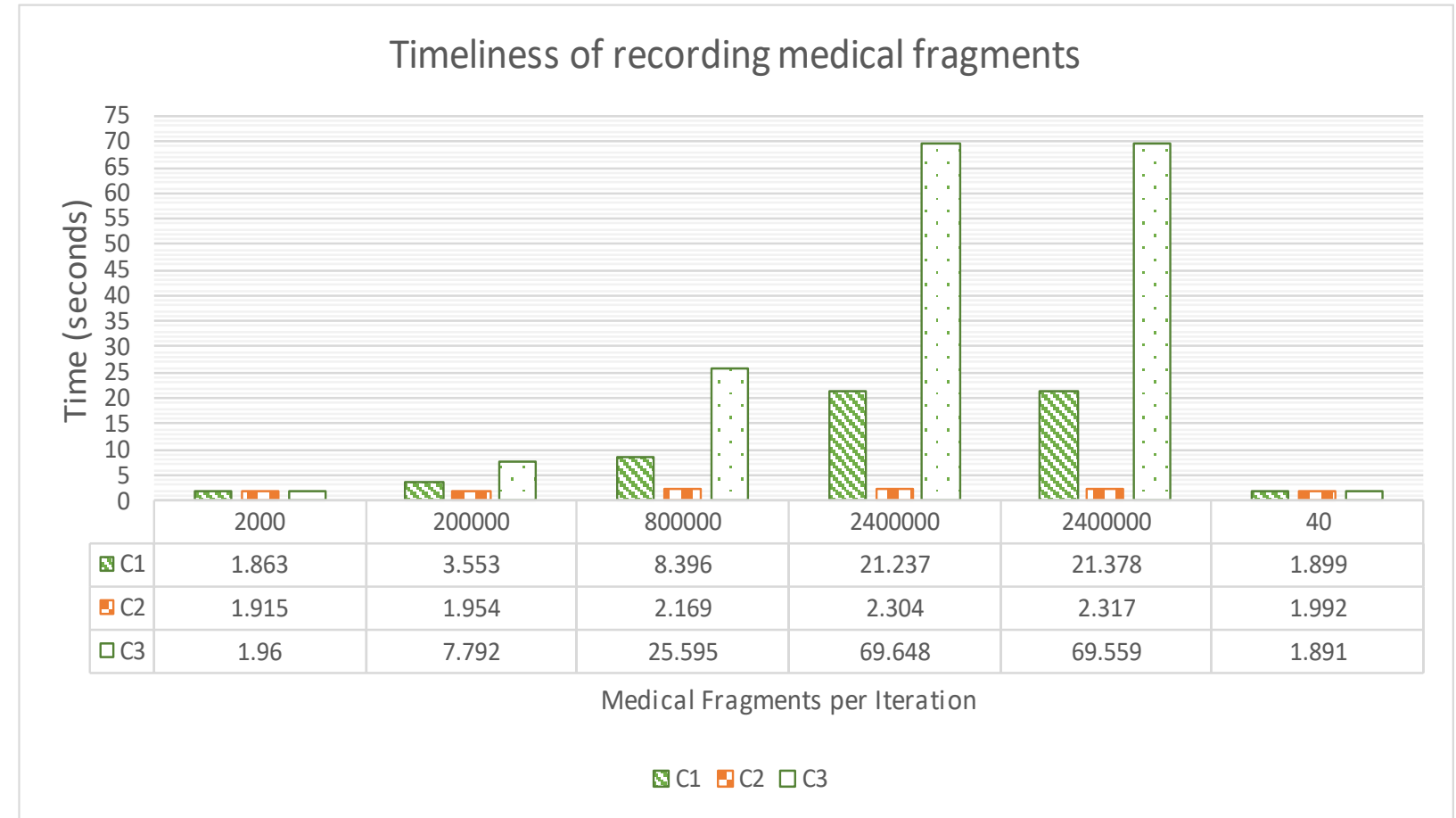| Criteria | Description | Metric |
|---|---|---|
| C1 | Time taken to insert medical fragment file into HDFS | Time |
| C2 | Time taken to insert medical fragment bridging information, linking Patient Id with fragment id into HDFS | Time |
| C3 | Time taken to insert patient index part of L-Store into HDFS | Time |
| C4 | Time taken to create table schema in Hive | Time |
| C5 | Time taken to create medical fragment bridging table schema in Hive | Time |
| C6 | Time taken to create patient index table schema in Hive | Time |
| C7 | Time taken to retrieve all fragment ids for 1 user | Time |
| C8 | Time taken to retrieve all medical fragments for 1 user | Time |

# Experimental Setup

**Solution 3:** Semantic Reconciliation-on-Read
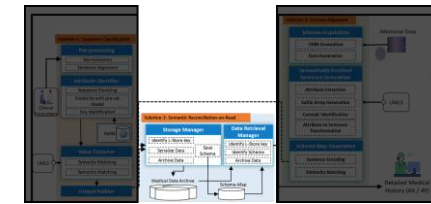
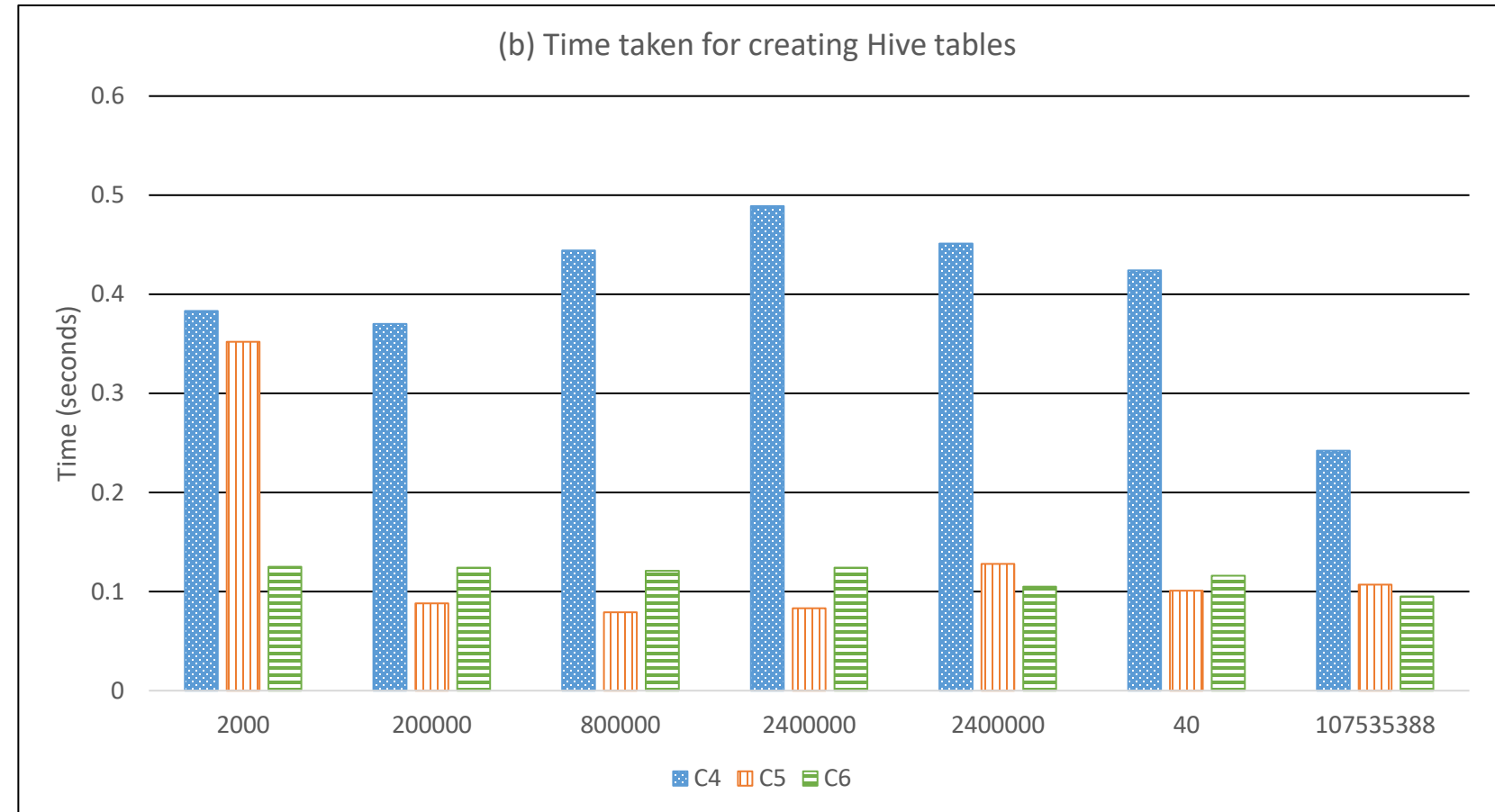| Criteria | Description |
|----------|-------------|
| C1 | Time taken to insert patient index part of L-Store into HDFS |
| C2 | Time taken to insert medical fragment bridging information, linking Patient Id with fragment id into HDFS |
| C3 | Time taken to insert medical fragment file into HDFS |

### Timeliness of recording medical fragments



| Medical Fragments per Iteration | 2000 | 200000 | 800000 | 2400000 | 2400000 | 40 |
|---------------------------------|------|--------|--------|---------|---------|-----|
| C1 | 1.863 | 3.553 | 8.396 | 21.237 | 21.378 | 1.899 |
| C2 | 1.915 | 1.954 | 2.169 | 2.304 | 2.317 | 1.992 |
| C3 | 1.96 | 7.792 | 25.595 | 69.648 | 69.559 | 1.891 |

# Experimental Setup

**Solution 3:** Semantic Reconciliation-on-Read
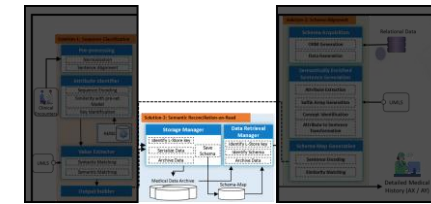
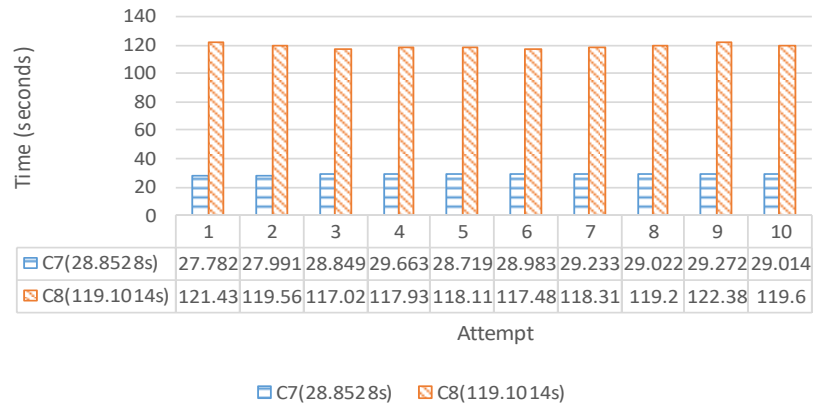| Criteria | Description |
|----------|-------------|
| C4 | Time taken to create table schema in Hive |
| C5 | Time taken to create medical fragment bridging table schema in Hive |
| C6 | Time taken to create patient index table schema in Hive |



(b) Time taken for creating Hive tables

# Experimental Setup

**Solution 3:** Semantic Reconciliation-on-Read
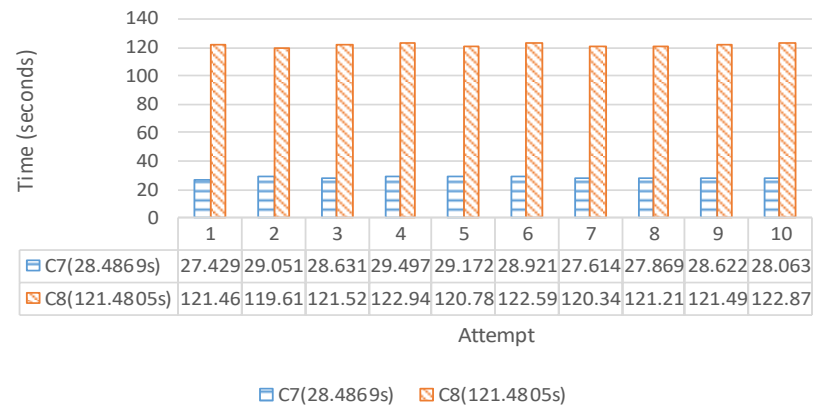
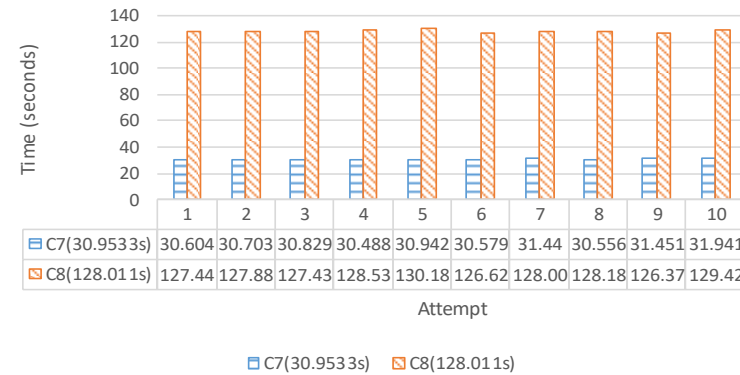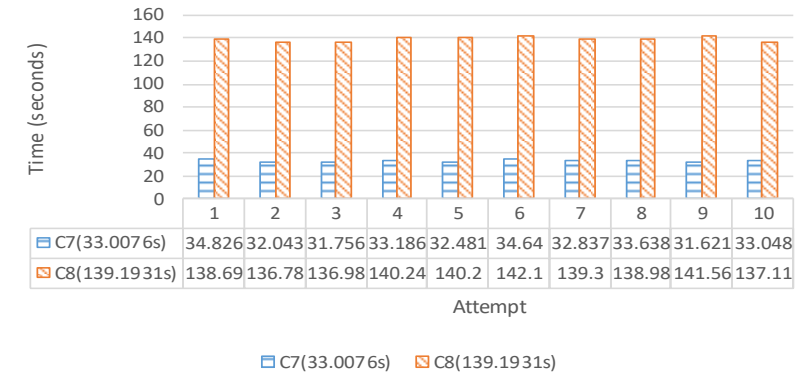| Criteria | Description | Metric |
|---|---|---|
| C7 | Time taken to retrieve all fragment ids for 1 user | Time |
| C8 | Time taken to retrieve all medical fragments for 1 user | Time |

## Iteration 1

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| C7(28.8528s) | 27.782 | 27.991 | 28.849 | 29.663 | 28.719 | 28.983 | 29.233 | 29.022 | 29.272 | 29.014 |
| C8(119.1014s) | 121.43 | 119.56 | 117.02 | 117.93 | 118.11 | 117.48 | 118.31 | 119.2 | 122.38 | 119.6 |

## Iteration 2

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| C7(28.4869s) | 27.429 | 29.051 | 28.631 | 29.497 | 29.172 | 28.921 | 27.614 | 27.869 | 28.622 | 28.063 |
| C8(121.4805s) | 121.46 | 119.61 | 121.52 | 122.94 | 120.78 | 122.59 | 120.34 | 121.21 | 121.49 | 122.87 |

## Iteration 3

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| C7(30.9533s) | 30.604 | 30.703 | 30.829 | 30.488 | 30.942 | 30.579 | 31.44 | 30.556 | 31.451 | 31.941 |
| C8(128.011s) | 127.44 | 127.88 | 127.43 | 128.53 | 130.18 | 126.62 | 128.00 | 128.18 | 126.37 | 129.42 |

## Iteration 4

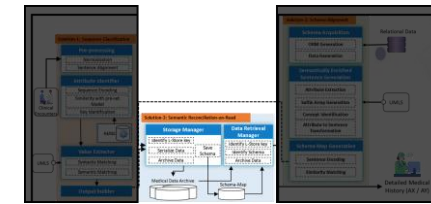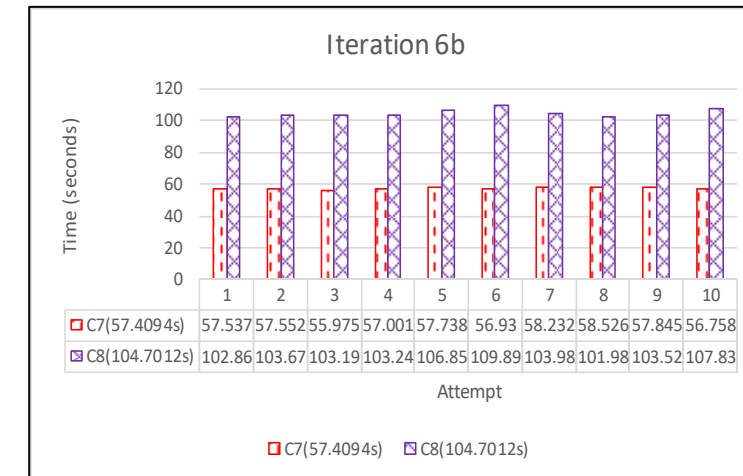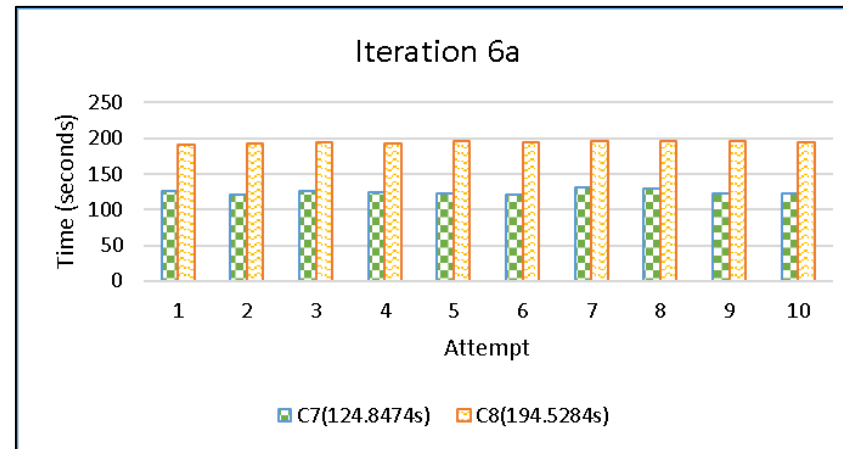| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| C7(33.0076s) | 34.826 | 32.043 | 31.756 | 33.186 | 32.481 | 34.64 | 32.837 | 33.638 | 31.621 | 33.048 |
| C8(139.1931s) | 138.69 | 136.78 | 136.98 | 140.24 | 140.2 | 142.1 | 139.3 | 138.98 | 141.56 | 137.11 |

# Experimental Setup

**Solution 3:** Semantic Reconciliation-on-Read
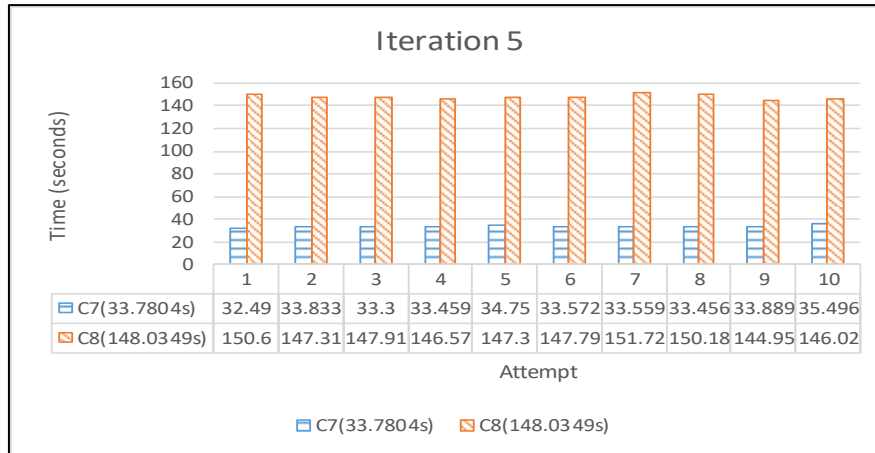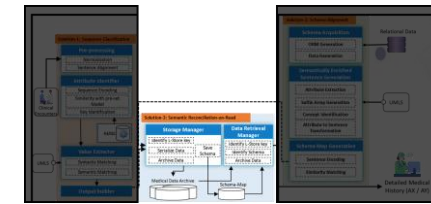
| Criteria | Description | Metric |
|---|---|---|
| C7 | Time taken to retrieve all fragment ids for 1 user | Time |
| C8 | Time taken to retrieve all medical fragments for 1 user | Time |

## Iteration 5

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| C7(33.7804s) | 32.49 | 33.833 | 33.3 | 33.459 | 34.75 | 33.572 | 33.559 | 33.456 | 33.889 | 35.496 |
| C8(148.0349s) | 150.6 | 147.31 | 147.91 | 146.57 | 147.3 | 147.79 | 151.72 | 150.18 | 144.95 | 146.02 |

## Iteration 6a

C7(124.8474s)   C8(194.5284s)

## Iteration 6b

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| C7(57.4094s) | 57.537 | 57.552 | 55.975 | 57.001 | 57.738 | 56.93 | 58.232 | 58.526 | 57.845 | 56.758 |
| C8(104.7012s) | 102.86 | 103.67 | 103.19 | 103.24 | 106.85 | 109.89 | 103.98 | 101.98 | 103.52 | 107.83 |

# Experimental Setup

**Solution 3:** Semantic Reconciliation-on-Read



| Criteria | Description | Metric |
|---|---|---|
| C7 | Time taken to retrieve all fragment ids for 1 user | Time |
| C8 | Time taken to retrieve all medical fragments for 1 user | Time |

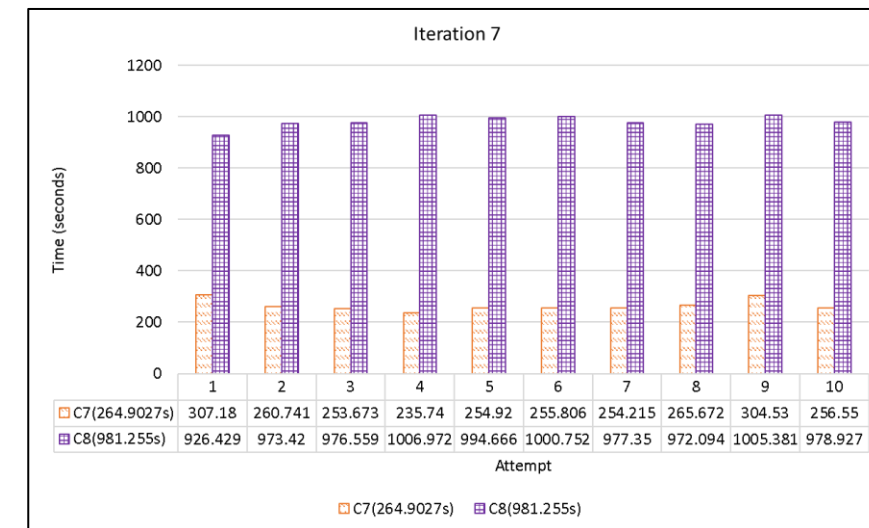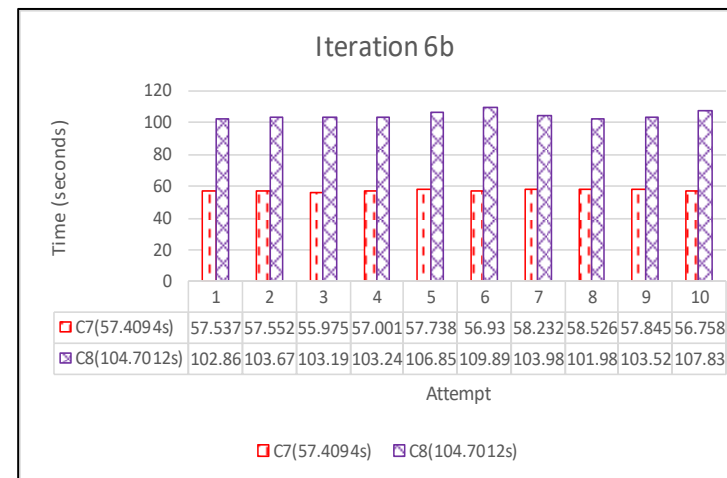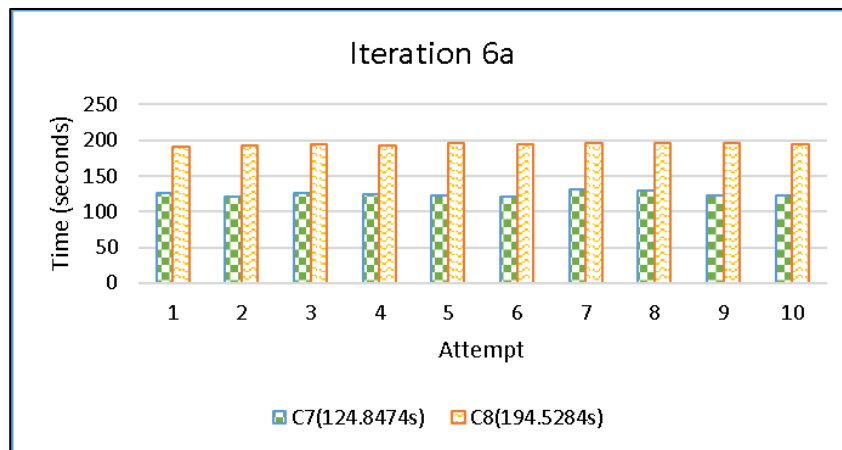# Solution 1: Paper Status

## Information Processing & Management
Supports *open access*

11 CiteScore | 7.466 Impact Factor

**em** Information Processing and Management

Fahad Ahmed Satti ⌄ | Logout

Home | Main Menu | Submit a Manuscript | About ⌄ | Help ⌄

← **Submissions Needing Revision for Author**

Click 'File Inventory' to download the source files for the manuscript. Click 'Revise Submission' to submit a revision of the manuscript. If you Decline To Revise the manuscript, it will be moved to the Declined Revisions folder.

IMPORTANT: If your revised files are not ready to be submitted, do not click the 'Revise Submission' link.

Page: 1 of 1 (1 total submissions)　　　Results per page  10 ⌄

| Action ➕ | 🏷✕ | Manuscript Number ▲ | Title ▲ | Initial Date Submitted ▼ | Date Revision Due ▲ | Status Date ▲ | Current Status ▲ | View Decision ▲ |
|---|---|---|---|---|---|---|---|---|
| Action Links | | IPM-D-22-01601R1 | A Semantic Sequence Similarity based approach for Extracting Medical Entities from Clinical Conversations | Jul 14, 2022 | Nov 02, 2022 | Oct 12, 2022 | Revise | Minor Revision |

Page: 1 of 1 (1 total submissions)　　　Results per page  10 ⌄