# Medical Language Models for Data Scientists

October 2025

# Building Patient Journeys & Cohorts

# Clinical Data is Spread in Multiple Sources

**Multi-Modal**

**Unstructured**

**Not Normalized**

**Not Consistent**

**Not Certain**

**Continuously Updating**

# Some Clinical Data is Only Available in Unstructured Text

## A Question-and-Answer System to Extract Data From Free-Text Oncological Pathology Reports

"Accuracies for predicting group-level site and histology codes were 93.5% and 97.6% respectively."

## Large language models to identify social determinants of health in electronic health records

"Our models identified 93.8% of patients with adverse SDoH, while ICD-10 codes captured 2.0%."

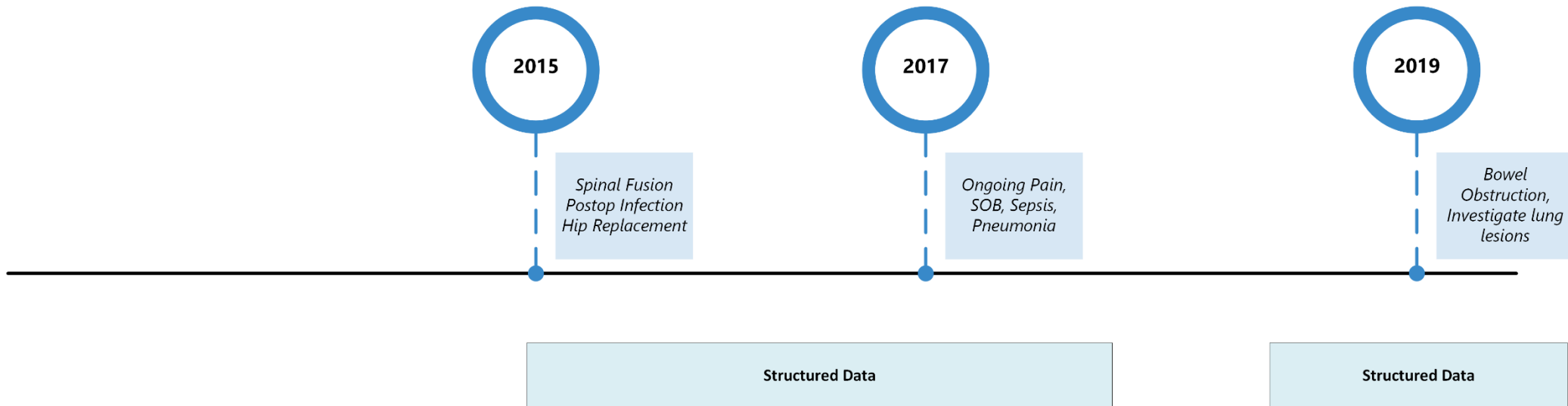## The Importance of Information Extraction from Unstructured Clinical Data in Pharmacoepidemiology

"the number of observed suicidality and self-harm events doubled with the addition of unstructured EHR data."
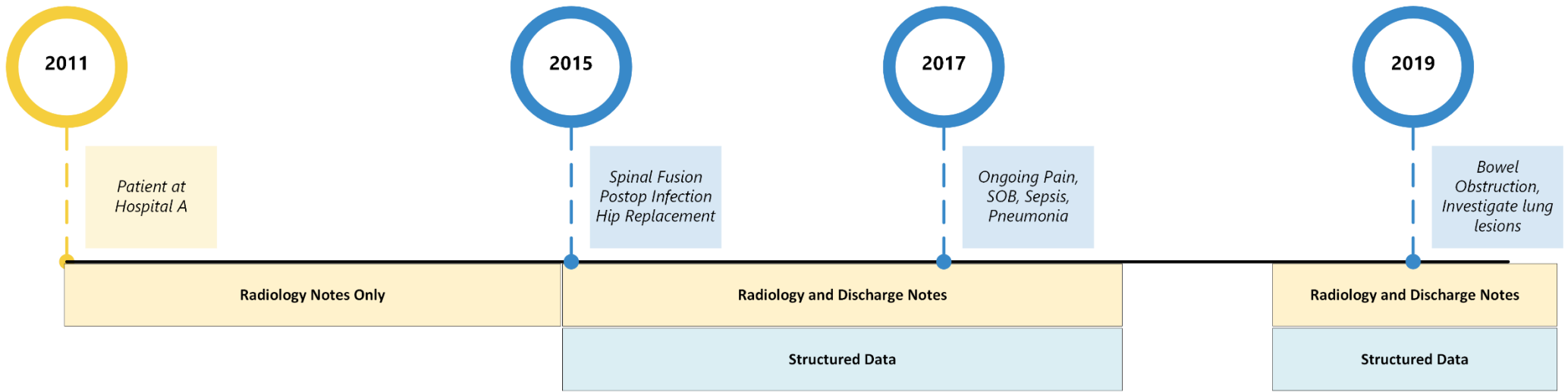
## An Assessment of Family History Information Captured in an Electronic Health Record

58.7% of the observations from the Neurology Admission Note contained family history, versus only 5.2% of structured records.
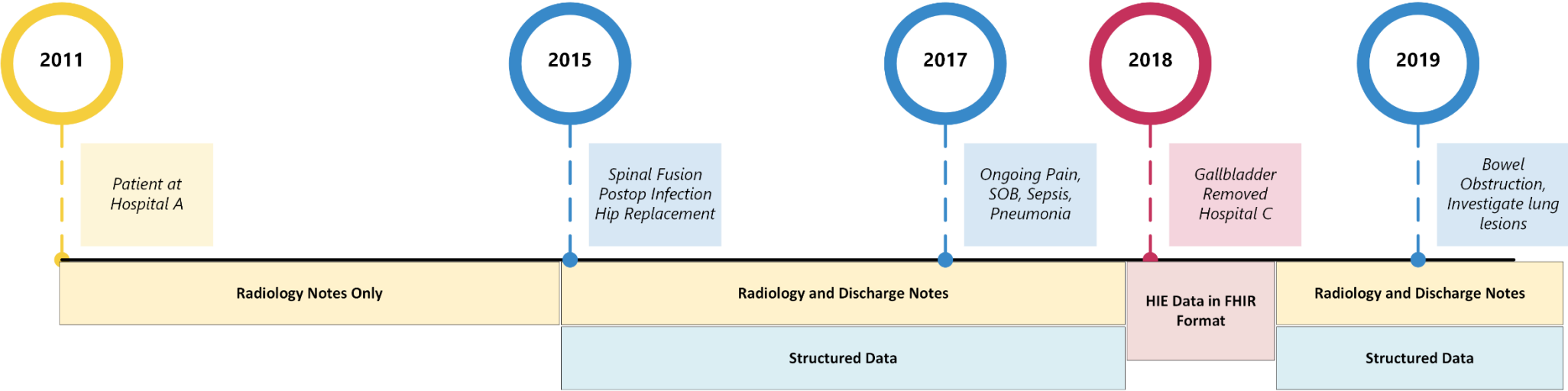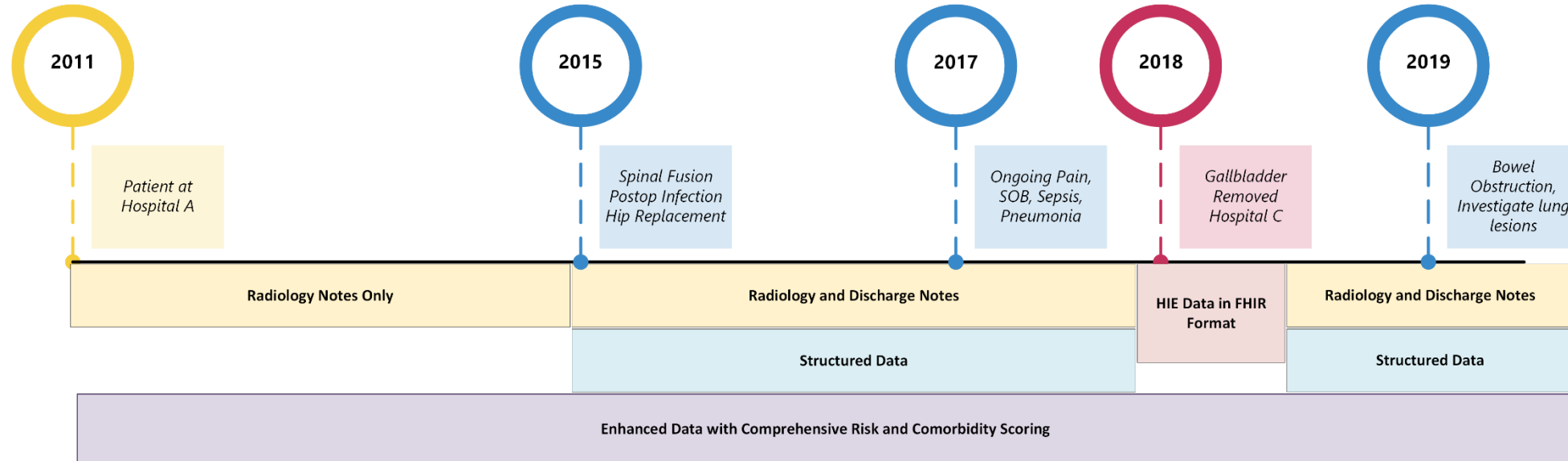
# One Patient's Journey: Structured EHR Data



**2015** — Spinal Fusion, Postop Infection, Hip Replacement

**2017** — Ongoing Pain, SOB, Sepsis, Pneumonia

**2019** — Bowel Obstruction, Investigate lung lesions

**Structured Data**

**Structured Data**

# Adding Radiology & Discharge Notes



**2011**

*Patient at Hospital A*

**2015**

*Spinal Fusion Postop Infection Hip Replacement*

**2017**

*Ongoing Pain, SOB, Sepsis, Pneumonia*

**2019**

*Bowel Obstruction, Investigate lung lesions*

**Radiology Notes Only**

**Radiology and Discharge Notes**

**Structured Data**

**Radiology and Discharge Notes**

**Structured Data**

# Adding FHIR Resources

# Current LLMs Can't Work on Complex Queries

**Q: "Find patients diagnosed with back pain that have had spinal fusion."**

- **RAG** can't find relevant information
- **Text2SQL** hallucinates in real-world DB settings, or build **queries that fail**
- Lack of **consistency**
- **Democratize** Cohort creation

```sql
WITH
-- Identify tuberculosis diagnosis concepts
tb_diagnosis_concepts AS (
    SELECT c.concept_id
    FROM concept c
    JOIN concept_ancestor ca ON ca.descendant_concept_id = c.concept_id
    WHERE ca.ancestor_concept_id IN (
        -- Add the root concept ID for tuberculosis and its descendants
        SELECT concept_id
        FROM concept
        WHERE concept_name = 'Tuberculosis Family'
        -- Ensure you replace 'Tuberculosis Family' with the correct name if different
    )
        AND c.standard_concept = 'S'
),

-- Identify tuberculosis drug concepts
tb_drug_concepts AS (
    SELECT c.concept_id
    FROM concept c
    JOIN concept_ancestor ca ON ca.descendant_concept_id = c.concept_id
    WHERE ca.ancestor_concept_id IN (
        -- Add the root concept ID for tuberculosis treatments and its descendants
        SELECT concept_id
        FROM concept
        WHERE concept_name = 'Tuberculosis Treatment'
        -- Ensure you replace 'Tuberculosis Treatment' with the correct name if
different
    )
        AND c.standard_concept = 'S'

), ● ● ●
```

# Our Approach

# Semantic Information Extraction (Healthcare NLP)



The patient is a 40-years-old black woman with [breast cancer | CANCER_DX | PRESENT]. She started [smoking | SMOKING_STATUS | PAST] when she was 20 years old, but she quit several years ago. Her mother died of [breast cancer | CANCER_DX | FAMILY] at age 55.

| begin | end | entity_type | assertion | confidence |
|-------|-----|-------------|-----------|------------|
| 47 | 59 | Cancer_Dx | Present | 0.9992 |
| 74 | 80 | Smoking_Status | Past | 0.9310 |
| 160 | 172 | Cancer_Dx | Family | 1.0000 |

# Terminology Server: Resolving to Standard Codes

```
{
  "url": "https://fhir/Mutation#assessed.gene",
  "valueCodeableConcept": {
    "coding": [
      {
        "system": "http://ncit.nci.nih.gov",
        "code": "C17757",
        "display": "EGFR"
      }
    ]
  }
},
{
  "url": "https://fhir/Mutation#assessed.referenceSeq",
  "valueCodeableConcept": {
    "coding": [
      {
        "system": "http://ncbi.nlm.nih.gov/CCDS",
        "code": "5514.1",
        "display": "CCDS 5514.1"
      }
    ]
  }
},
{
  "url": "https://fhir/Mutation#assessed.variant",
  "valueCodeableConcept": {
    "coding": [
      {
        "system": "http://www.hgvs.org/mutnomen",
        "code": "c.2369C>T",
        "display": "c.2369C>T"
      },
      {
        "system": "http://www.hgvs.org/mutnomen",
        "code": "p.T790M",
        "display": "T790M"
      }
```

| ner_chunk | entity | snomed_code | resolved_text |
|---|---|---|---|
| Catheterization of left heart | Procedure | 67629009 | Catheterization of left heart |
| selective coronary angiogram | Test | 33367005 | Coronary angiography |
| common femoral angiogram | Test | 4701000087107 | Angiography of right femoral artery |
| StarClose closure of right common femoral artery | Procedure | 310621009 | Patch repair of femoral artery |

# Merging & Deduplicating Facts to Build a Patient Graph

# Making Clinical Inferences and Calculations

## Rule-based Medical Calculation

**Patient Note**

A 68-year-old man with the left hemiparesis from 2 h previously visited the emergency room. His medical history included hypertension and bilateral emphysema due to heavy smoking. Vital sign assessment revealed tachycardia; examination of the heart revealed atrial [...]

**Question**

What is the patient's **CHA2DS2-VASc score**?

**Explanation**

The patient is 68 years old. Because the age is between 65 and 74, one point added to the score, making the current total 0 + 1 = 1. The patient's gender is male so no points are added to the current total, keeping the total at 1. The patient history for congestive heart [...]

**Final Answer**

7

## Equation-based Medical Calculation

**Patient Note**

The patient was a 20-year-old previously healthy woman. She was a university student. Her height and body weight were 168.1 cm and 52.2 kg, respectively. She ingested bamboo salt (about 150 grams ) in a day for the purpose of digestion and weight reduction [...]

**Question**

What is the patient's **albumin corrected anion gap** in mEq/L?

**Explanation**

The formula for computing a patient's albumin corrected anion gap is: anion gap (in mEq/L) + 2.5 * (4 - albumin (in g/dL)). The formula for computing a patient's anion gap is: sodium (mEq/L) - (chloride (mEq/L)+ bicarbonate (mEq/L)). The concentration of sodium [...]

**Final Answer**

19.25

**MedCalc-Bench: Evaluating Large Language Models for Medical Calculations**

Nikhil Khandekar, Qiao Jin, Guangzhi Xiong, Soren Dunn, Serina S Applebaum, Zain Anwar, Maame Sarfo-Gyamfi, Conrad W Safranek, Abid A Anwar, Andrew Zhang, Aidan Gilson, Maxwell B Singer, Amisha Dave, Andrew Taylor, Aidong Zhang, Qingyu Chen, Zhiyong Lu

# The OMOP CDM

**Observational Medical Outcomes Partnership - Common Data Model**

Enhancing Healthcare through Data, since 2009

**Foundation:** Part of the Observational Health Data Sciences and Informatics (OHDSI) initiative.

**Objective:** Utilize open-source data solutions to improve human health via large-scale analysis.

**Purpose:** Standardize the structure and content of observational healthcare data.

**Methods:**

- Through pseudonymisation and common data quality assessments, the OMOP-CDM provides a robust framework for converting complex EMR data into a standardised format.

- By securely sharing de-identified and aggregated data and conducting analyses across multiple OMOP-converted databases, patient-level data is securely firewalled within its respective local site.

Seamless EMR data access: Integrated governance, digital health and the OMOP-CDM, Feb 2024.
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10882353/



**OHDSI** OBSERVATIONAL HEALTH DATA SCIENCES AND INFORMATICS

*"an international collaborative whose goal is to create and apply open-source data analytic solutions to a large network of health databases to improve human health and wellbeing"*

# The OMOP CDM

**O**bservational **M**edical **O**utcomes **P**artnership - **C**ommon **D**ata **M**odel



*CDM v5.4*

- 39 tables
- 433 fields
- 7 categories

https://ohdsi.github.io/CommonDataModel/cdm54.html

# Data Ingestion

# Cohort Building from Prompts

What would you like to know?

Show patients who were diagnosed with back pain and who have had spinal fusion ⌄

| | condition_name | condition_source | procedure_name |
|---|---|---|---|
| 0 | Backache | EHR billing record | Lumbar and lumbosacral fusion of the posterior column |
| 1 | Backache | EHR billing record | Lumbar and lumbosacral fusion of the posterior column |
| 2 | Backache | EHR billing record | Fusion of 2 or more Lumbar Vertebral Joints with Autologo |
| 3 | Backache | EHR billing record | Fusion of 2 or more Lumbar Vertebral Joints with Synthetic |
| 4 | Backache | EHR billing record | Fusion of Lumbosacral Joint with Autologous Tissue Subst |
| 5 | Backache | 14260897-DS-5 | Arthrodesis |
| 6 | Backache | 14260897-DS-5 | Lumbar and lumbosacral fusion of the posterior column |
| 7 | Low back pain | EHR billing record | Fusion of Lumbar Vertebral Joint with Interbody Fusion De |
| 8 | Low back pain | EHR billing record | Lumbar and lumbosacral fusion by anterior technique |
| 9 | Low back pain | EHR billing record | Lumbar and lumbosacral fusion of the posterior column |

# Behind the scenes: Multi-agent system

Find patients diagnosed with `back pain` that have had `spinal fusion`

Concept resolver: Find concept id for given entity.

**Back pain**: Condition (SNOMED 194133)
**Spinal fusion**: Procedure (SNOMED 4177164)

Build query for OMOP CDM

Retrieve records and make reply

# Patient Level Analysis: Chat With Your Data

# Track Provenance of Information



- EHR System
- FHIR document
- Raw clinical notes
- Other

# OMOP CDM (PostgreSQL)

Records for one patient

| table_name 🔒 name | row_count 🔒 bigint |
|---|---:|
| note_nlp | 5852 |
| observation | 211 |
| visit_occurrence | 151 |
| note | 151 |
| condition_occurrence | 138 |
| measurement | 76 |
| person | 11 |
| procedure_occurrence | 11 |
| drug_exposure | 8 |