

Towards Patient-Driven Phenotyping and Similarity for Precision Medicine

Tiffany J. Callahan
Computational Biosciences Program
University of Colorado Anschutz Medical Campus

Olivier Bodenreider
Lister Hill National Center for Biomedical Communications
National Library of Medicine



Electronic Medical Records

▶ Digital version of a patient's medical history:

- ▶ Inpatient notes
- ▶ Labs and physical exams
- ▶ Prescribed medications
- ▶ Diagnoses and procedures
- ▶ Treatment plans
- ▶ Discharge instructions

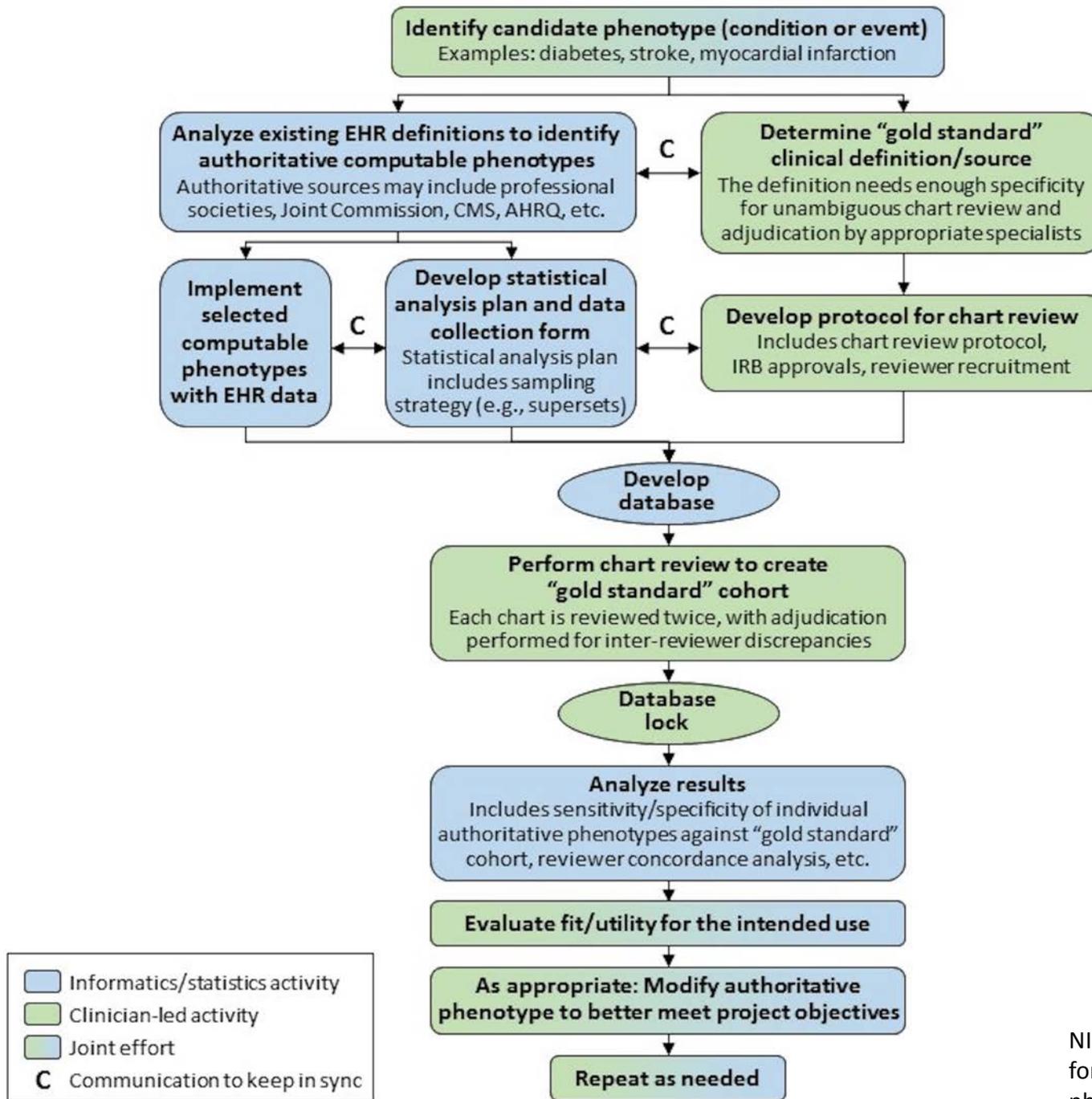
The screenshot displays the Allscripts Professional EHR interface for a patient named Deanna Daley. The top navigation bar includes the patient's name, date of birth (2/15/1981), gender (29y Female), insurance (Blue Cross & Blue Shield), and ID number (#210). Below this, a summary of patient status is provided: Status: Active, Usual: Neuron, Nate, Ref: Amblin, Arthur B. MD, Allergies: Latex, No Known Drug Allergies, Guarantor: Deanna Daley, Web Account, Marital Status: Single, and Blood Type: O+ (Patient reported). The main content area is divided into several panels: 'Face Sheet' on the left, 'Medical History: Newest to oldest' in the center, 'Encounters: By Type, Newest to Oldest' on the right, and 'Medications: All, Newest to Oldest' at the bottom right. The 'Medical History' panel shows a list of conditions including 'MIGRAINE WITH AURA, NON-INTRACTABLE (346.00)', 'COMMON MIGRAINE WITHOUT MENTION OF INTRACTABLE MIGRAINE', 'Allergy: Latex: Rash, Hives', 'No Known Drug Allergies', 'Family: Negative Family History of: CVA, TIA, Temporal Arteritis, Major Depression', 'Social: No Drug Use, Non Smoker/No Tobacco Use, Caffeine Use: 2-3 cups coffee / day, Alcohol Use: Occasional alcohol use', 'Pregnancy/Birth: Pregnancies (Gravida) [11/2006]: Gravida 1', 'Past Surgical: Appendectomy [1999]', and 'Hospitalizations - Dates/Reasons: 1996 - appendectomy, 2003 - child birth'. The 'Encounters' panel lists 'Consultation', 'Office Visit', '[Open Encounter]', 'Chart Attachments', 'Labs/Procedures', and 'Referral Letter'. The 'Medications' panel shows 'Current Medications' including 'Maxalt 5MG, 1 (one) Tablet at onset of headache. May repeat in 2 hours. M', 'Amoxicillin 125MG/5ML, 1 For Suspension daily, 1 For Suspension, 1 day s', and 'Cipro 250MG, 1 Tab BID, 14 Tab, 07/05/2010, No Refill. Active.', along with 'Administered Medications' and 'Previous Medications' including 'Imitrex 5MG/ACT, 2 (two) Not Specified q2h PRN, 30 days starting 08/12/2010'. The 'Orders' panel shows 'Future' orders such as '10/6/2010: METABOLIC PANEL, COMPREHENSIVE (80053) [Future Order]', '10/6/2010: SED RATE ERYTHROCYTE (85651) [Future Order]', '10/6/2010: TSH (THYROID STIMULATING HORMONE) (84443) [Future Order]', and '10/6/2010: CBC, PLATELETS & AUT DIFF (85025) [Future Order]'. The right sidebar contains 'Actions' like 'Menu', 'Send Message', 'Launch', 'Print', and 'Queues' with various counts for 'Received Charts', 'Appointments', 'Open Encounters', 'Result Notifications', 'Messages', 'Web Messages', 'Refill Requests', and 'eRefill Requests'.

The big promise of lies in large-scale use, automatically feeding clinical research, quality improvement, and clinical phenotyping.

Clinical Phenotyping

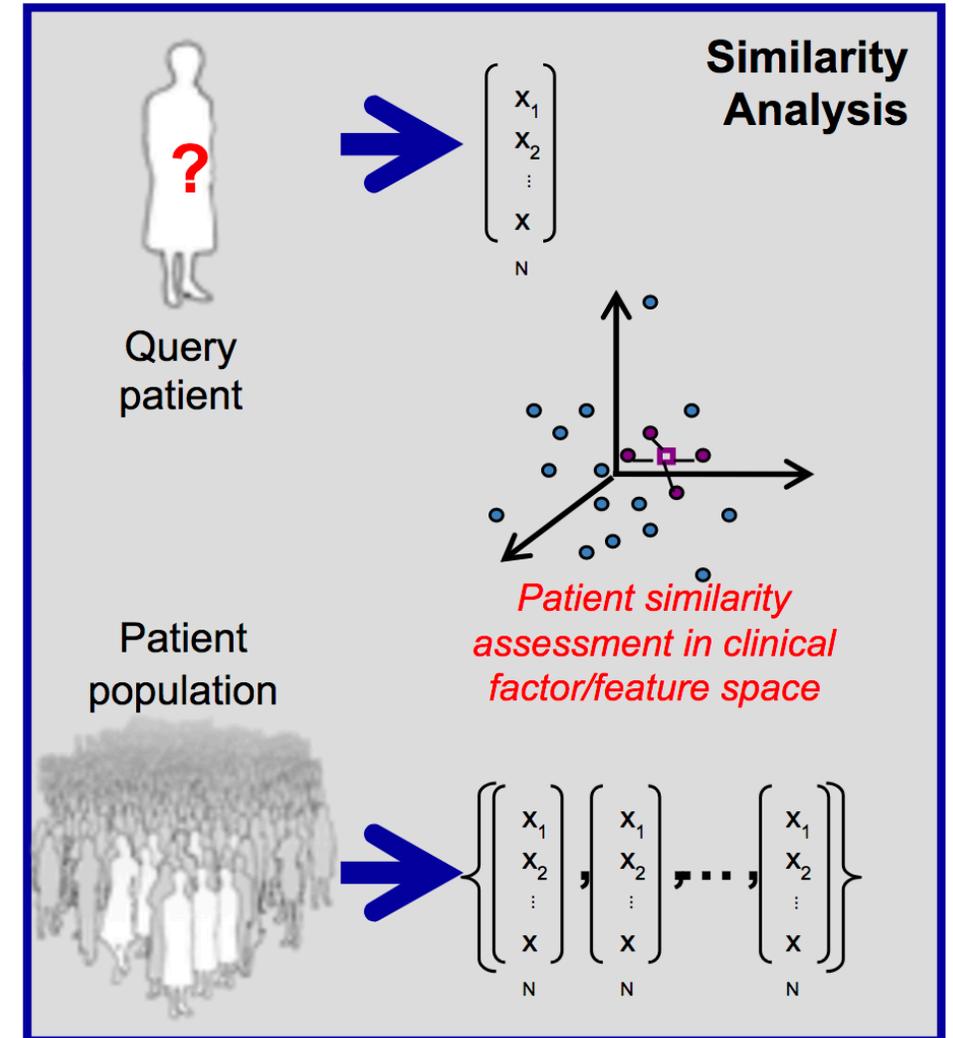
- ▶ **Goal:** identify cohorts of patients with specific clinical features characteristic of a disease of interest
- ▶ Typical approaches:
 - ▶ Rule-based
 - + interpretable, fast to implement, good results on limited datasets
 - requires expert knowledge and multiple iterations, not easily generalizable
 - ▶ Natural language processing and text mining
 - + rich data not found in other sources
 - sensitive to misspelling/bad grammar, redundancy, ambiguity
 - ▶ Machine learning
 - + many standardized approaches, easy implementation, robust
 - curse of dimensionality; difficult with rare disease/small patient cohorts

PheKB



Motivation

- Traditional approaches are good at providing information on the “average patient”
- What evidence can physicians use when trying to treat a patient whose symptoms deviate from average?
- **Patient Similarity:** derive insights from patients that are similar to an index patient to provide personalized predictions¹
- Diagnostic cohort identification
 - Drug repurposing
 - Identify and tailor treatment recommendations



Approach

1. Similarity function

- ▶ Data-driven; automatic
- ▶ Pediatric data - OMOP CDM v5

2. Clustering

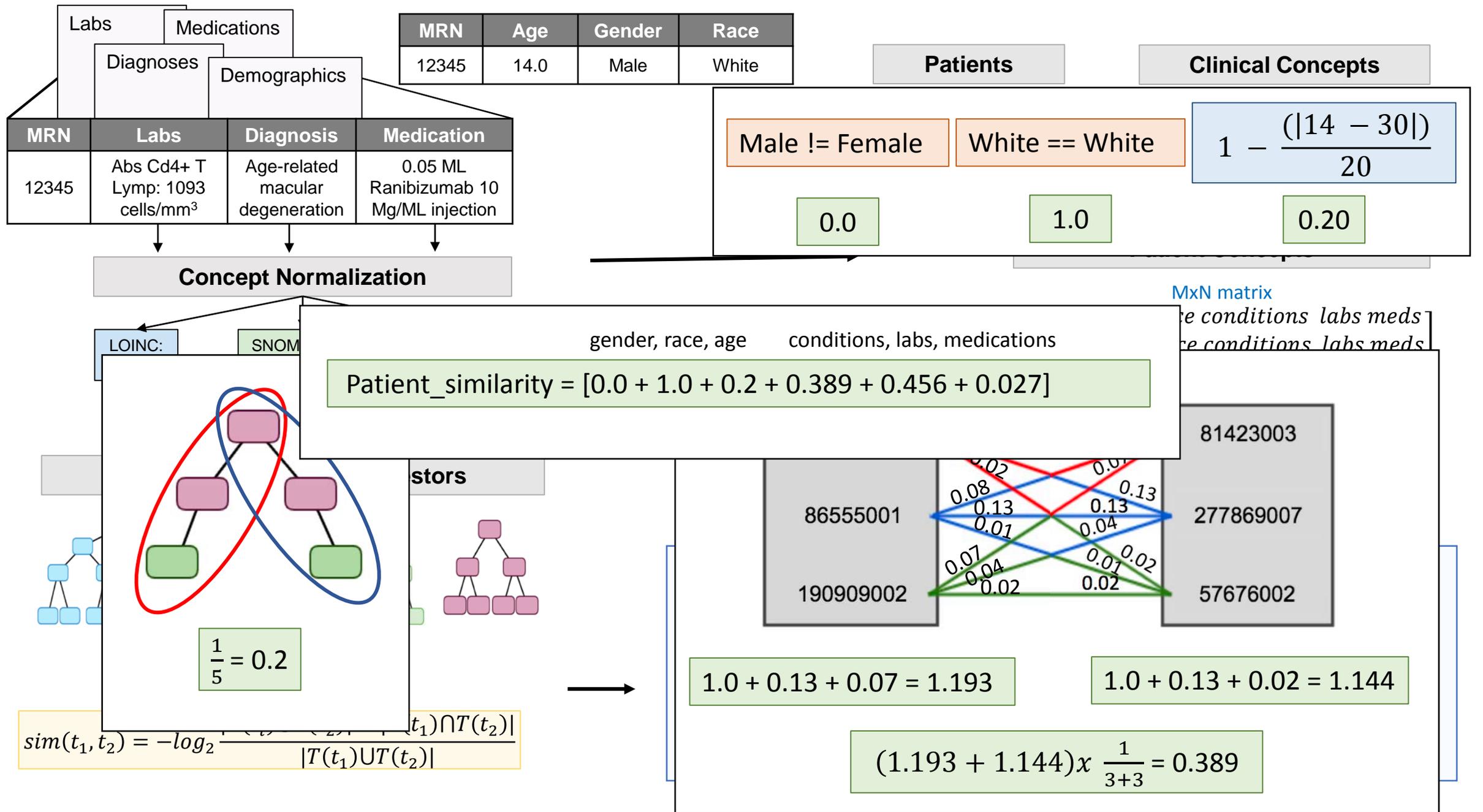
- ▶ Similarity function-driven

3. Cluster identification/labeling

- ▶ Clinical terminologies/value sets
- ▶ Biomedical Knowledgebase
- ▶ Literature

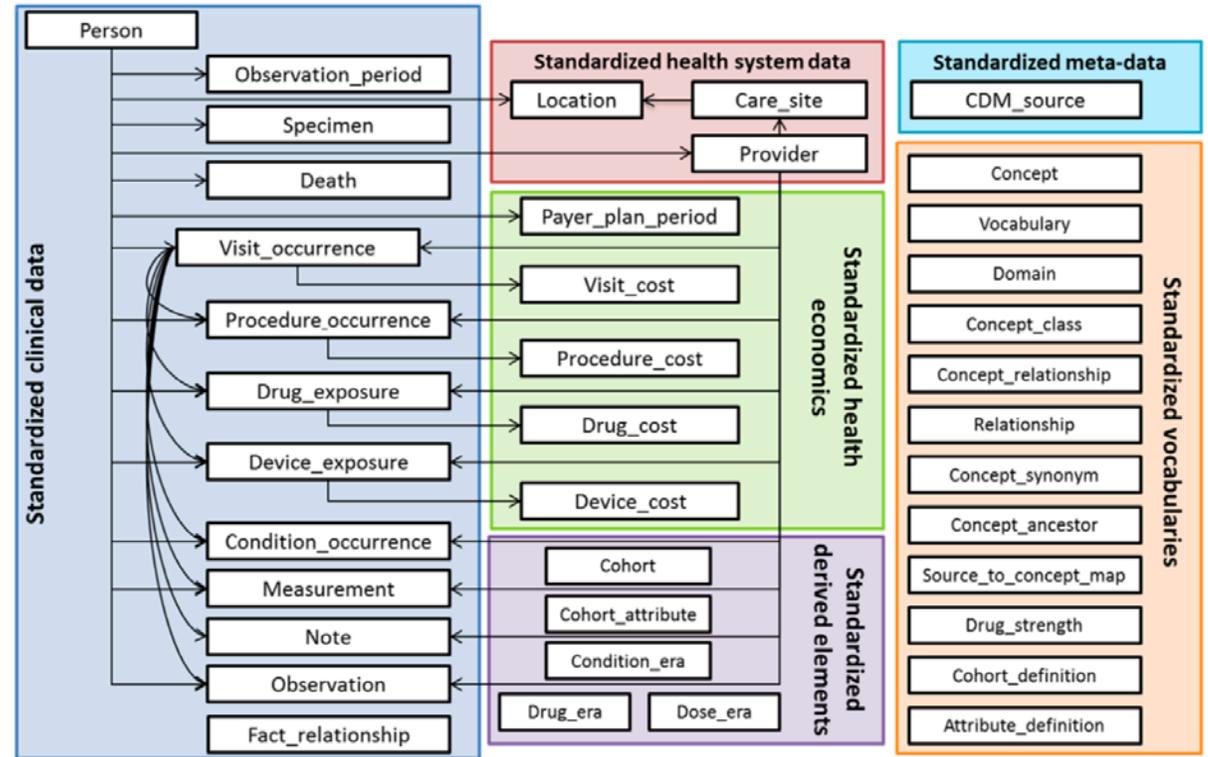
4. Evaluation

- ▶ Compare to PheKB clusters
- ▶ Verify algorithm reproducibility across data warehouses



Evaluation

- ▶ Children's Hospital of Colorado EHR data
 - ▶ De-identified (COMIRB # 15-0445)
- ▶ PEDSnet OMOP version 5
 - ▶ Concepts normalized to standardized terminologies
- ▶ Test Case – 2 groups (N = 20)
 - ▶ Huntington's Chorea (ICD-9-CM 333.4)
 - ▶ Cystic Fibrosis (ICD-9-CM 722.0)



```

SELECT
  person_id,
  condition_source_value,
  COUNT(condition_source_value) AS count
FROM omop5deid.condition_occurrence
WHERE condition_source_value LIKE '333.4 %'
GROUP BY person_id, condition_source_value
ORDER BY count DESC
LIMIT 10
;
    
```

Row	person_id	condition_source_value	count
1	151130	333.4 Huntington's chorea	125
2	322136	333.4 Huntington's chorea	47
3	456778	333.4 Huntington's chorea	21
4	65447	333.4 Huntington's chorea	19
5	53666	333.4 Huntington's chorea	16
6	434388	333.4 Huntington's chorea	15

Huntington's Chorea

- ▶ Fatal disorder caused by breakdown of nerve cells in the brain
- ▶ 30,000 Americans have been diagnosed
- ▶ Symptoms include:
 - ▶ Personality, mood changes
 - ▶ Unsteadiness, poor coordination
- ▶ Diagnoses ($\bar{x} = 320.2$; unique = 334)
- ▶ Laboratory tests ($\bar{x} = 114.1$, unique = 119)
- ▶ Medications ($\bar{x} = 1528.7$, unique = 177)

[Hdsa.org/what-is-hd/](https://hdsa.org/what-is-hd/)

Row	co_person_id	OMOP_concept	SNOMED_code	source_code	count
1	151130	374341	58756001	333.4 Huntington's chorea	125
2	151130	435642	29164008	784.5 Other speech disturbance	59
3	151130	441277	25766007	315.32 Mixed receptive-expressive language disorder	57
4	151130	441417	281016006	781.3 Lack of coordination	40
5	151130	440086	192127007	314.00 Attention deficit disorder without mention of hyperactivity	38
6	151130	440374	191736004	300.3 Obsessive-compulsive disorders	22
7	151130	442077	197480006	300.00 Anxiety state, unspecified	17
8	151130	436077	248290002	315.9 Unspecified delay in development	17
9	151130	438409	406506008	314.9 Unspecified hyperkinetic syndrome	16
10	151130	436073	69322001	298.9 Unspecified psychosis	16
11	151130	379782	5158005	307.23 Tourette's disorder	14

Cystic Fibrosis

- ▶ Genetic disease that causes mucus buildup resulting in persistent lung infection and difficulty breathing
- ▶ >30,000 people diagnosed worldwide
- ▶ Symptoms include:
 - ▶ Coughing, wheezing, frequent lung infections
 - ▶ Poor growth, male infertility
- ▶ Diagnoses ($\bar{x} = 982.5$, unique = 447)
- ▶ Laboratory tests ($\bar{x} = 3104.4$, unique = 124)
- ▶ Medications ($\bar{x} = 3120.3$, unique = 392)

<https://www.cff.org/What-is-CF/About-Cystic-Fibrosis/>

Row	co_person_id	OMOP_concept	SNOMED_code	source_code	count
1	388727	434615	81423003	277.00 Cystic fibrosis without mention of meconium ileus	538
2	388727	254320	86555001	277.02 Cystic fibrosis with pulmonary manifestations	406
3	388727	378253	25064002	784.0 Headache	352
4	388727	442077	197480006	300.00 Anxiety state, unspecified	343
5	388727	194325	190909002	277.03 Cystic fibrosis with gastrointestinal manifestations	272
6	388727	436096	82423001	338.29 Other chronic pain	215
7	388727	4174281	277869007	031.0 Pulmonary diseases due to other mycobacteria	194
8	388727	377889	15188001	389.9 Unspecified hearing loss	182
9	388727	77074	57676002	719.40 Pain in joint, site unspecified	150
10	388727	317009	195967001	493.90 Asthma, unspecified type, unspecified	111
11	388727	4032376	108305003	V57.1 Care involving other physical therapy	95

Conclusions

- ▶ Developed a patient similarity algorithm
 - ▶ Data-driven
 - ▶ Composite semantic similarity for heterogeneous data types
 - ▶ Adjust weights to customize by use case
 - ▶ Promising initial proof of concept with pediatric EHR is promising
- ▶ Limitations
 - ▶ Small test group, need to scale to larger groups
 - ▶ Limited evaluation
 - ▶ Several unmapped Generic Product Identifiers
- ▶ Future Work
 - ▶ Explore alternative semantic similarity algorithms
 - ▶ Optimize algorithm
 - ▶ Develop machine learning approach to determine patient similarity attribute weights

Acknowledgments

- ▶ Dr. Olivier Bodenreider
- ▶ Drs. Paul Fontelo and Clement McDonald
- ▶ Summer Follows: Ann Cirincione and Raja Cholan
- ▶ Dr. Michael Kahn and Health Data Compass Team