

### Ontologizing Health Systems at Scale: Making Translational Discovery a Reality

Tiffany J. Callahan MPH, PhD Candidate

OHDSI Community Call - September 2020



### 98% of Hospitals Use Electronic Health Records

"...now is the time to create smarter healthcare systems in which the best treatment decisions are computationally learned from electronic health record data by deep-learning methodologies"<sup>1</sup>



### Creating Smarter Healthcare Systems

#### Current research use of medical record data enables

- Automatic triage of medical conditions
- Prediction of pertinent patient risk factors
- Identification of emerging or important pathogens

### Creating Smarter Healthcare Systems

**Current research use of medical record data enables** 

- Automatic triage of medical conditions
- Prediction of pertinent patient risk factors
- Identification of emerging or important pathogens

#### Connecting medical record data with biomedical knowledge could enable

- Integration of patient -omics data
- Identification of causal biomarkers
- Mechanistic explanations for specific patient characteristics





# **BARRIER 1**

## Medical Records are not Connected to Biomedical Knowledge

### Knowledge Graphs



### Knowledge Graphs Challenges

### **Design Challenges**

- Multiple approaches to modeling biomedical knowledge<sup>1-3</sup>
- Decisions impact downstream learning

#### **Implementation Challenges**

- Methods have varying functional, logical, and semantic consequences<sup>1</sup>
- Semantic Web standard is expressive, but unwieldy and of uncertain benefit<sup>2</sup>

**PheKnowLator:** A Python framework to build large-scale heterogeneous biomedical knowledge graphs



### Human Disease Mechanisms

**Ontologies** 12 OBO ontologies 366,846 classes 3,923,625 axioms

#### **Edge Data** 22 Linked Open Data

2 Experimental 34 edge-types

Validation PhD Molecular Biologist



## Summary

- First fully custom biomedical knowledge graph construction framework
- Evaluation in progress
- Promising applications
  - Biomedical hypergraphs  $\rightarrow$  Joslyn, Aksoy, Callahan et al., arXiv. 2020
- Full testing and continuous integration; Dockerized

# **BARRIER 2**

## Medical Records were not Built to Facilitate Translational Research

### Clinical Research Challenges

#### Problem

- Hospital databases/EHRs incomplete and not standardized<sup>1-3</sup>
- Diagnosis codes used for billing, not research<sup>2,4</sup>

### Solution

- CDMs help different data sources speak the same language
- Ontologies contain meaningful representations of molecular knowledge

**EHRs** = Electronic Health Records; **CDMs** = Common Data Models

### Ontologies

**Ontology**  $\rightarrow$  Graph of entities and relationships<sup>1,2</sup>

- Domain-specific
- Community consensus
- Hierarchical

### **Open Biomedical Ontologies**

- Phenotypes (Human Phenotype Ontology)
- Diseases (Disease Ontology)
- Chemicals/Metabolites/Hormones (ChEBI)
- Anatomical Entities (UBERON)
- Cell Types (Cell Line Ontology)



## **Ontologies Translate Clinical Concepts**



## Existing Mapping Work



- Manually curated 2,923 LOINC lab tests to HPO
- Validated with 15,681 patients with respiratory complaints  $\rightarrow$  known asthma biomarkers
  - Abnormal metabolism/vitamin metabolism
  - ↑ Red blood cell count and VLDL cholesterol concentration

## **Mapping Strategies**

"Chronic deep venous thrombosis of right calf"

1:1 Mapping Approach:

"Chronic deep venous thrombosis of right calf"  $\rightarrow$  "Abnormality of the calf"

#### 1:Many Mapping Approach:

"Chronic deep venous thrombosis of right calf"  $\rightarrow$  "Chronic"

"Deep venous thrombosis" "Abnormality of the calf" "Right"

## Objective

- Extend and expand existing mapping work  $\rightarrow$  more than 1:1 mappings
- Map ontologies to OMOP
- Hospital scale, disease-agnostic

**OMOP2OBO:** the first health system-wide integration and alignment between OMOP standardized clinical terminologies and eight OBO biomedical ontologies





## Available Data

Clinical Domain	OMOP Table	Concept Class	Concept Vocabularies	Ontologies
Conditions	Condition Occurrence	Conditions	SNOMED CT Source Codes: "Mapped to", "Mapped From", "Concept poss_eq from" "Concept same_as from"	Human Phenotype Ontology Mondo Disease Ontology
Medications	Drug Exposure Immunizations	Drugs Ingredients	RxNorm Standard Non-Standard	ChEBI Protein Ontology NCBITaxon Vaccine Ontology
Measurements	Measurements	Measurements	LOINC Standard Non-Standard	Human Phenotype Ontology ChEBI Uber Anatomy Ontology Protein Ontology NCBITaxon Cell Ontology

## Available Data

Clinical Data - OMOP CDM					
Concept ID*	138994				
Source Code	snomed:109995007 icd10cm:D46.9 mesh:D009190 icd9cm:238.75				
Concept Name	Myelodysplastic syndrome				
Concept Synonym	Myelodysplastic syndrome (disorder)   Myelodysplastic syndrome (clinical)				

\*OMOP Data also includes all ancestors

#### **Ontology Data - Open Biomedical Ontologies**

Class ID	HP_0002863
Class Label	Myelodysplasia
Class Synonym - Broad - Exact - Narrow - Related	Hypoplastic myelodysplasia
Class Definition	Clonal hematopoietic stem cell disorders characterized by dysplasia (ineffective production) in one or more hematopoietic cell lineages, leading to anemia and cytopenia
Class DbXRef	UMLS:C1851971 MSH:D009190 SNOMEDCT_US:109995007

## Mapping Approach

#### **Mapping Strategies**

- Database Cross-References
- Exact String Match
  - Labels
  - Synonyms
- Bag-of-Words + TF-IDF Weighting
  - Labels
  - Synonyms
  - Definitions (ontologies only)

## Mapping Categories

Category	Definition
Automatic Exact - Concept	Exact label or synonym, dbXRef, or expert validated mapping @ concept-level; 1:1
Automatic Exact - Ancestor	Exact label or synonym, dbXRef, or expert validated mapping @ concept ancestor-level; 1:1
Manual Exact - Concept Similarity	Concept similarity score suggested mapping manually verified; 1:1
Automatic Constructor - Concept	Exact label or synonym, dbXRef, cosine similarity, or expert validated mapping @ concept-level; 1:Many
Automatic Constructor - Ancestor	Exact label or synonym, dbXRef, cosine similarity, or expert validated mapping @ concept-level; 1:Many
Manual Constructor	Hand mapping created using expert suggested resources; 1:Many
Manual	Hand mapping created using expert suggested resources; 1:1
UnMapped	No suitable mapping or not mapped type

## Mapping Categories

Category	Definition
Automatic Exact - Concept	Exact label or synonym, dbXRef, or expert validated mapping @ concept-level; 1:1
Automatic Exact - Ancestor	Exact label or synonym, dbXRef, or expert validated mapping @ concept ancestor-level; 1:1
Manual Exact - Concept Similarity	Concept similarity score suggested mapping manually verified; 1:1
Automatic Constructor - Concept	Exact label or synonym, dbXRef, cosine similarity, or expert validated mapping @ concept-level; 1:Many
Automatic Constructor - Ancestor	Exact label or synonym, dbXRef, cosine similarity, or expert validated mapping @ concept-level; 1:Many
Manual Constructor	Hand mapping created using expert suggested resources; 1:Many
Manual	Hand mapping created using expert suggested resources; 1:1
UnMapped	No suitable mapping or not mapped type

Example: Apraxia (OMOP\_132342)

#### Mapping Type: Automatic Exact - Concept

Ontology Concept: Apraxia (HP\_0002186)

#### Mapping Evidence:

CONCEPT\_DBXREF:snomed\_68345001 OBO\_LABEL-OMOP\_CONCEPT\_LABEL:apraxia OBO\_LABEL-OMOP\_CONCEPT\_SYNONYM:apraxia CONCEPT\_SIMILARITY:HP\_0002186\_1.0 ANCESTOR\_DBXREF:snomed\_68345001 OBO\_LABEL-OMOP\_ANCESTOR\_LABEL:apraxia

### **Conditions Occurrence Mappings**

Маррі	ng Type	Condition Concept ID	Phenotype (Human Phenotype Ontology)	Disease (Mondo Disease Ontology)	
Automatic	Concept Macular Hole (0MOP_4338894)		Macular Hole (HP_0011508)	Macular Hole (DOID_7633)	
Mapping	Ancestor	Major histocompatibility complex class I deficiency (OMOP_4100979)	Macular Hole (нр_0011508) Severe Combined Immunodeficiency (нр_0004430) Abnormality of the large intestin morphology (нр_0005210)	Severe Combined Immunodeficiency (DOID_627)	
Manual	Similarity	Malakoplakia of colon (OMOP_4024250)	Abnormality of the large intestine morphology (HP_0005210)	Colonic Disease (DOID_5353)	
Mapping	Constructor	Macular edema and retinopathy due to type 2 diabetes mellitus (OMOP_45770830)	AND   Type II Diabetes Mellitus (HP_0005978)   Macular Edema (HP_0040049)   Retinopathy (HP_0000488)	AND   Type 2 Diabetes (DOID_9352)   Macular Retinal Edema (DOID_4449)	

. . . . . . . . . . . . . . . . . .

### **Conditions Occurrence Mappings**

- UMLS CUIs (MRCONSO 2020AB) and Semantic Types (MRSTY 2020AB)
- 28,129 unique condition occurrence codes

**Evaluation:** 20% of manually mapped concepts verified by clinicians; several iterations

- 24,285 OMOP concepts  $\rightarrow$  4,661 Phenotypes
- 19,664 OMOP concepts  $\rightarrow$  3,614 Diseases

Ontology	Automatic Exact - Concept	Automatic Exact - Ancestor	Manual Exact - Concept Similarity	Automatic Constructor - Concept	Automatic Constructor - Ancestor	Manual Constructor	Manual
HPO	3465	2851	1055	174	1825	9477	5438
MONDO	4965	6103	484	723	2301	3109	1979

### Drug Exposure Mappings

**Mapping:** Ingredients and mechanism of action (DrugBank and CHEMBL)

Example: balsalazide (OMOP\_934262)

Ingredient mapping:

- Automatic Exact Concept
- balsalazide (CHEBI\_267413)

#### Mechanism of Action:

- Automatic Constructor Concept
- **agonist:** peroxisome proliferator-activated receptor gamma (PR\_P37231)
- **Inhibitor:** prostaglandin G/H synthase 2 (PR\_P35354), prostaglandin G/H synthase 1 (PR\_P23219), arachidonate 5-lipoxygenase-activating protein (PR\_P20292)
- organism: homo sapien (NCBITaxon\_9606)

### Drug Exposure Mappings

• 11,937 drug-ingredient combinations  $\rightarrow$  1,697 unique ingredients

Evaluation: 20% of manually mapped concepts verified by clinical pharmacist; 3 iterations

- 1,618 OMOP concepts  $\rightarrow$  1,422 Chemicals, hormones, or metabolites
- 139 OMOP concepts  $\rightarrow$  91 Proteins
- 317 OMOP concepts  $\rightarrow$  39 Organisms
- 127 OMOP concepts  $\rightarrow$  54 Vaccines/Immunizations

Ontology	Automatic Exact - Concept	Automatic Exact - Ancestor	Manual Exact - Concept Similarity	Automatic Constructor - Concept	Automatic Constructor - Ancestor	Manual Constructor	Manual	Unmapped
CHEBI	1176	3	64	9	1	58	306	80
PRO	2	0	10	0	0	8	119	1558
NCBITaxon	20	0	0	0	0	10	287	1380
VO	92	0	8	1	0	1	25	1570

### Measurement Mappings

- Leverage LOINC scale and results type
- Data-drive confirmation

Somatotropin [Mass/volume] in Serum or Plasma (OMOP\_3023709)

- OR blood serum(UBERON\_0001977) | blood plasma (UBERON\_0001969)
- Somatotropin (PR\_000007968)
  - Normal: NOT Abnormality of circulating hormone level (HP\_0003117)
  - Low: Growth Hormone Deficiency (HP\_0000824)
  - High: Growth Hormone Excess (HP\_0000845)

Transitional cells [Presence] in Urine sediment by Light microscopy (OMOP\_3028475)

- Urine (UBERON\_0001088), Transitional epithelial cell (CL\_0000244)
  - Negative: NOT Increased urinary transitional epithelial cell count (HP\_0032214)
  - Positive: Increased urinary transitional epithelial cell count (HP\_0032214)

### Measurement Mappings

• 4,382 lab tests, 11,072 lab test results

#### **Evaluation:**

- 270 results verified by 3 MDs, 1 epidemiologist
- 15% verified by a biocurator, 3 iterations
  - 10,888 OMOP concepts  $\rightarrow$  920 Phenotypes | 10,876 OMOP concepts  $\rightarrow$  25 Anatomical entities
  - 1,075 OMOP concepts  $\rightarrow$  27 Cell Types | 9,710 OMOP concepts  $\rightarrow$  338 Chemicals, Hormones
  - 8,269 OMOP concepts  $\rightarrow$  194 Organisms | 4,842 OMOP concepts  $\rightarrow$  Proteins

Ontology	Automatic Exact - Concept	Automatic Exact - Ancestor	Manual Exact - Concept Similarity	Automatic Constructor - Concept	Automatic Constructor - Ancestor	Manual Constructor	Manual
HPO	6891	22	171	4	0	49	3751
UBERON	0	1612	0	0	0	3307	5957
CL	0	227	207	0	87	30	524
CHEBI	5	3535	1098	151	241	922	3758
NCBITaxon	0	693	0	0	0	459	7117
PRO	0	44	259	0	0	175	4364

## OMOP2OBO

• Determine coverage of mappings on two independent samples

Data Source	Domain	Unique Concepts	Coverage
	Conditions	5608	92.13%
	Medications	Unique Concepts   5608   4426   15056   23972	96.36%
	Conditions	15056	79.78%
UCHEalth	Medications	23972	91.35%

## Summary

- First hospital scale mapping between OMOP and the OBOs
- Preliminary coverage examined on adult and ICU populations
- Added over 200 new concepts to HPO
- Adopted by National COVID-19 Cohort (N3C)



### Next Steps:

- Coverage study in subset of Concept Prevalence data
- Comparing to Juan Banda's Mappings

### ACKNOWLEDGEMENTS



University of Colorado Anschutz Medical Campus

#### Trivisory





Dr. Michael G. Kahn Professor of Clinical Informatics Co-Director of Colorado CTSI Director, Health Data Compass Assoc Director, Colorado CPM

Dr. Lawrence E. Hunter Professor of Pharmacology, Computer Science (CU Boulder), Director of the Computational

Bioscience Program



Dr. Tellen D. Bennett Assoc Professor of Pediatric Critical Care Co-Director, Analytics Core, CU Data Science to Patient Value Assoc Director, Informatics Core, CCTSI



Adrianne L. Stefanski, PhD CU Anschutz



Peter Robinson, MD

Professor

Jackson Labs





Nicole Vasilevsky, PhD Research Assistant Professor OHSU

**Pharmacology Experts** 

Xingmin Aaron Zhang, PhD Post Doctoral Fellow Jackson Labs

#### **Clinical Experts**



James Feinstein, MD Associate Professor General Pediatrics CU Anschutz



Blake Martin, MD Jordan Wyrwa, DO Pediatric Critical Care Physical Medicine and Medicine Fellow Rehabilitation Resident CU Anschutz CU Anschutz



Emily Swenson CU Anschutz Medical Student



Kelsey Andrews CU Anschutz Medical Student



Katy Trinkley, PharmD Associate Professor CU Anschutz

Catherine Derington, PharmD Post Doctoral Fellow CU Anschutz

Jessica Sinclair, PharmD Pharmacy Resident Purdue

Organizations











Work funded by the NIH, No. T15LM009451

#### Translational Research Experts