



# Standardizing Electronic Health Records to Improve Medical Research and Analysis

## The OMOP CDM

---

DATE: AUGUST 22, 2019    PRESENTED BY: ALYSSA HUQUE AND AARON BERGER, INTERNS



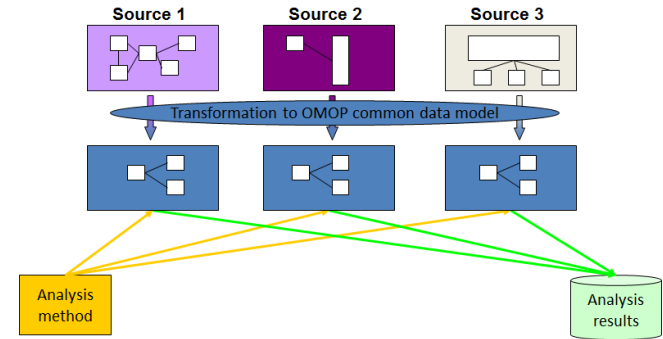
# Background

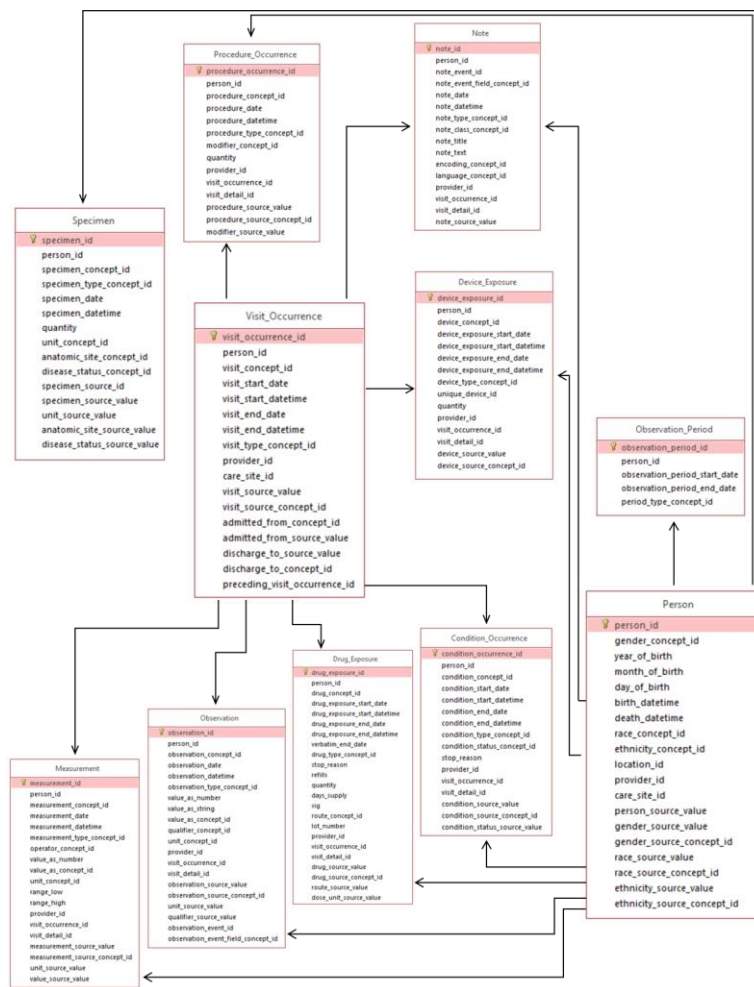
- Passage of the Health and Information Technology for Economic and Clinical Health Act (HITECH Act) of American Recovery and Reinvestment Act of 2009 (ARRA – Obama’s economic stimulus bill) incentivized medical institutions to digitize health records
- Inconsistencies of data captured
- Desire to standardize EHRs
  - Secondary research
  - Collaboration and compatibility with other data sets
  - Increased accuracy, reproducibility



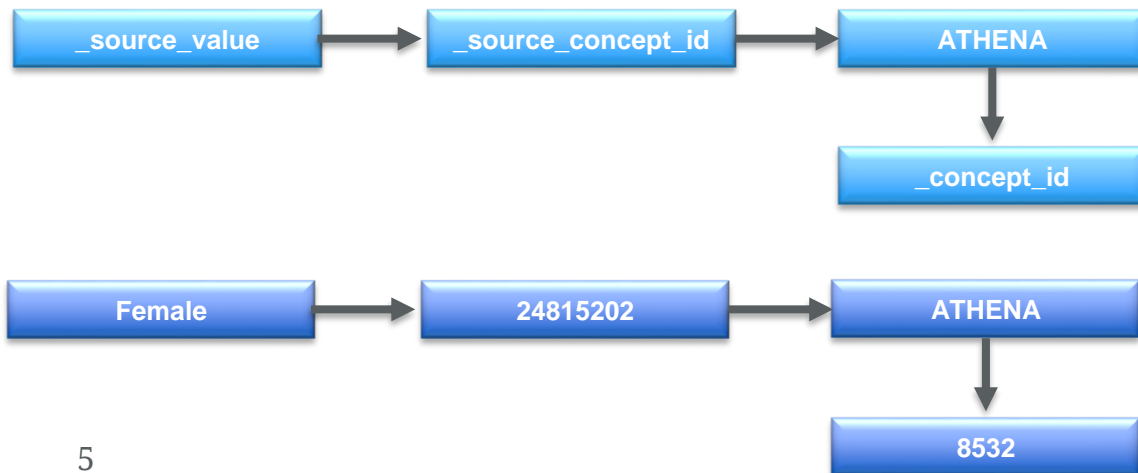
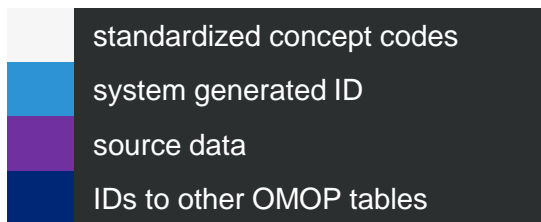
# OHDSI's OMOP CDM

- Standardizes data
  - Common format
  - Common representation
- Perform systematic and standard analytics
- Support research to identify and evaluate associations
- Allows data analysis from multiple, disparate data sources
  - Handles administrative claims and EHR





# Example OMOP Person Table



Field	Type
person_id	integer
gender_concept_id	integer
year_of_birth	integer
month_of_birth	integer
day_of_birth	integer
birth_datetime	datetime
death_datetime	datetime
race_concept_id	integer
ethnicity_concept_id	integer
location_id	integer
provider_id	integer
care_site_id	integer
person_source_value	varchar(50)
gender_source_value	varchar(50)
gender_source_concept_id	integer
race_source_value	varchar(50)
race_source_concept_id	integer
ethnicity_source_value	varchar(50)
ethnicity_source_concept_id	integer

# Example OHSU Demographics Table

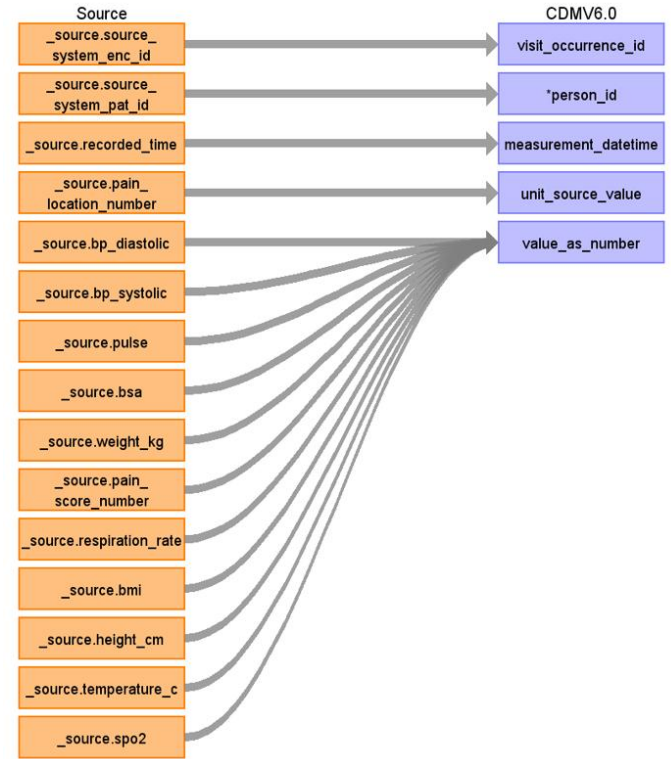
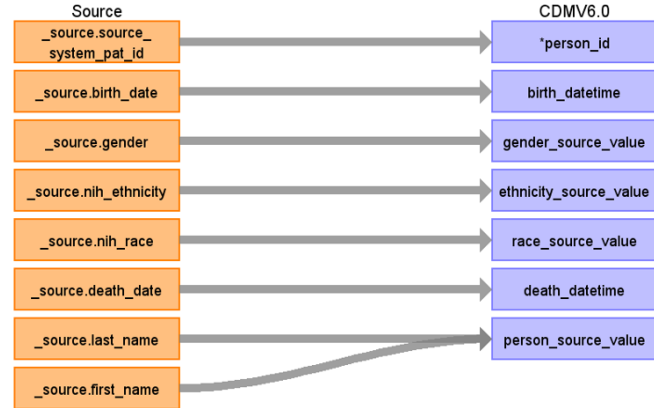
Field	Type
OHSU_MRN	integer
SOURCE_SYSTEM_PAT_ID	integer
CURRENT_AGE_YRS	integer
BIRTH_DATE	datetime
GENDER	varchar(50)
NIH_ETHNICITY	varchar(50)
NIH_RACE	varchar(50)
PATIENT_ALIVE	varchar(50)
DEATH_DATE	datetime
LAST_NAME	varchar(50)
FIRST_NAME	varchar(50)
ADDRESS_CITY	varchar(50)
ADDRESS_LINE1	varchar(50)
ADDRESS_LINE2	varchar(50)
ADDRESS_STATE	varchar(50)
ADDRESS_ZIP	integer
ADDRESS_PHONE	integer
EMAIL	varchar(50)
ADDRESS_COUNTY	varchar(50)
CURRENT_PCP	varchar(50)
BIO_SAMPL_OPT_OUT_FLG	varchar(50)
GENETC_OPT_OUT_FLG	varchar(50)





# Mapping Data

## Demographics vs. Vitals



# Issues with Mapping Data

- Loss of uniqueness of certain OHSU tables
- Fields can map but terminology won't

encounter_diagnosis			
BILLING_DX_FLAG	Y/N		
DX_DATE	YYYY-MM-DD	(CONDITION_OCCURRENCE) condition_start_datetime	datetime
DX_ICD	code		
DX_ICD10	code	(CONDITION_OCCURRENCE) condition_source_concept_id	varchar(50)
DX_ICD10_NAME	varchar, blank	(CONDITION_OCCURRENCE) condition_source_value	varchar(50)
DX_NAME	varchar		
ENC_DX_FLAG	Y/N		
FOLLOWUP_DX_FLAG	Y/N		
HOSP_ADMIT_FLAG	Y/N		
MED_HX_DX_FLAG	Y/N		
OHSU_MRN	integer		
ORDER_MED_DX_FLAG	Y/N		
ORDER_PROC_DX_FLAG	Y/N		
REFERRAL_DX_FLAG	Y/N		
SOURCE_SYSTEM_ENC_ID	integer	(CONDITION_OCCURRENCE) visit_occurrence_id	varchar(50)
SOURCE_SYSTEM_PAT_ID	Z + integer	(CONDITION_OCCURRENCE) person_id	integer





# Success in Mapping Data

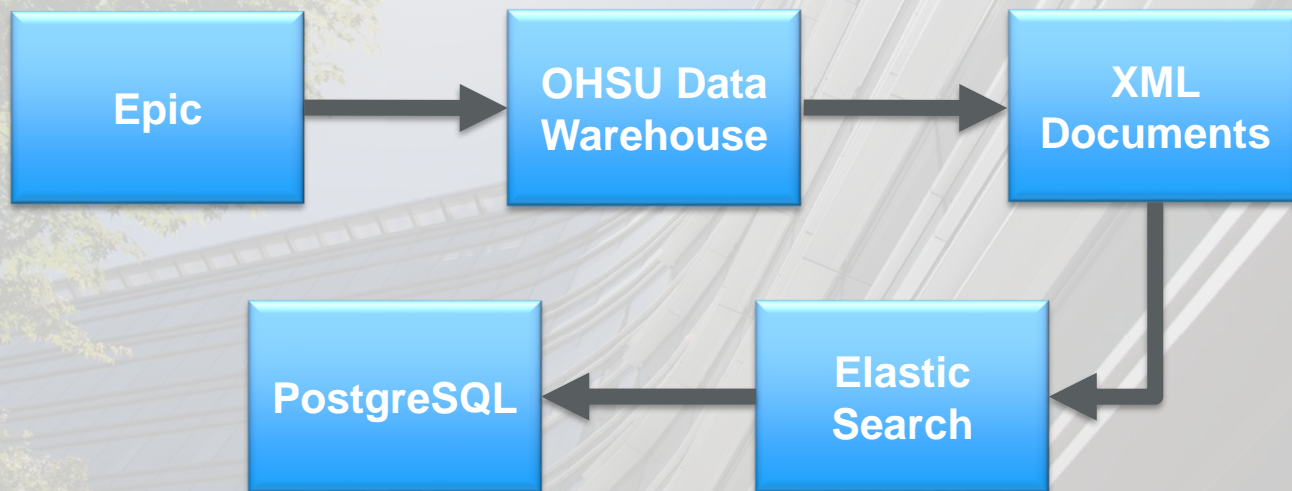
- 74% (198/268) source fields mapped successfully
- 62% (96/154) OMOP fields of pure data mapped successfully
  - Adding `_concept_ids` will increase percentages of mapped OMOP fields

# Aaron - PostgreSQL



- Use PostgreSQL to create CDM tables and source tables
- Extract data from Elasticsearch and import into created source tables
- Transform source data format into CDM data format
- Load newly formatted data into CDM tables
- Implement primary, foreign keys for relationality

# The Path of Aaron's Data



# PostgreSQL Tables and Fields

-- Fields that map from 'Demographics' to 'person'

```
INSERT INTO public.person(person_id, person_source_value, birth_datetime, death_datetime, gender_source_value, gender_concept_id, race_source_value, race_concept_id)
SELECT "SOURCE_SYSTEM_PAT_ID", "FULL_NAME", "BIRTH_DATE", "DEATH_DATE", "GENDER", "NIH_RACE",
FROM public."Demographics";
```

--Fields that map from 'Demographics' to 'location' (Health System Data Table)

```
INSERT INTO public.location(address_1, address_2, city, state, zip, county)
SELECT "ADDRESS_LINE1", "ADDRESS_LINE2", "ADDRESS_CITY", "ADDRESS_STATE", "ADDRESS_ZIP", "ADDRESS_COUNTY",
FROM public."Demographics";
```

--Fields that map from 'Notes' to 'note' (Clinical Data Table)

```
INSERT INTO public.note(note_datetime, note_text, note_id, person_id)
SELECT "NOTE_DATE", "NOTE_TEXT", "SOURCE_SYSTEM_NOTE_CSN_ID", "SOURCE_SYSTEM_PAT_ID"
FROM source."Notes";
```

--Fields that map from 'Problem\_List' to 'condition\_occurrence' (Clinical Data Table)

```
INSERT INTO public.condition_occurrence(person_id, condition_start_datetime, condition_end_datetime, condition_source_value, condition_concept_id)
SELECT "SOURCE_SYSTEM_PAT_ID", "DX_START_DATE", "DX_END_DATE", "DX_ICD10_NAME", "DX_ICD10"
FROM source."Problem_List";
```

	birth_datetime timestamp without time zone	death_datetime timestamp without time zone	gender_source_value character varying (50)	gender_concept_id integer	race_source_value character varying (50)	race_concept_id integer
1	1955-12-25 00:00:00	9999-12-31 00:00:00	FEMALE	[null]	WHITE	[null]
2	1956-12-25 00:00:00	9999-12-31 00:00:00	FEMALE	[null]	WHITE	[null]
3	1957-12-25 00:00:00	9999-12-31 00:00:00	FEMALE	[null]	WHITE	[null]
4	2002-12-25 00:00:00	9999-12-31 00:00:00	MALE	[null]	WHITE	[null]
5	1960-12-25 00:00:00	9999-12-31 00:00:00	FEMALE	[null]	WHITE	[null]
6	1943-12-25 00:00:00	9999-12-31 00:00:00	FEMALE	[null]	WHITE	[null]
7	1963-12-25 00:00:00	9999-12-31 00:00:00	FEMALE	[null]	WHITE	[null]
8	1989-12-25 00:00:00	9999-12-31 00:00:00	FEMALE	[null]	WHITE	[null]



# The Path of Alyssa's Data



# Alyssa – Python

- Goal: to develop and implement a software program that maps OHSU's EHR data stored in XML files to new XML files that are formatted and organized by the standards set by the OMOP CDM.



# Example OHSU EHR XML Document Demographics

```
<?xml version="1.0"?>
- <main>
  - <DATA_RECORD>
    <MEDICAL_RECORD_NUMBER>00000000</MEDICAL_RECORD_NUMBER>
    <PATIENT_ID>Z0000000</PATIENT_ID>
    <CURRENT_AGE_IN_YEARS>20</CURRENT_AGE_IN_YEARS>
    <BIRTH_DATE>04/21/1999</BIRTH_DATE>
    <GENDER>FEMALE</GENDER>
    <ETHNICITY>NOT HISPANIC OR LATINO</ETHNICITY>
    <RACE>ASIAN</RACE>
    <ALIVE>Y</ALIVE>
    <DEATH_DATE>12/31/1999</DEATH_DATE>
    <LAST_NAME>HUQUE</LAST_NAME>
    <FIRST_NAME>ALYSSA</FIRST_NAME>
    <CITY>WASHINGTON</CITY>
    <ADDRESS_LINE1>1600 PENNSYLVANIA AVENUE</ADDRESS_LINE1>
    <ADDRESS_LINE2>N/A</ADDRESS_LINE2>
    <STATE>DC</STATE>
    <ZIP>20500</ZIP>
    <PHONE>2024561414</PHONE>
    <EMAIL>ALYHUQUE&#64GMAIL.COM</EMAIL>
    <COUNTY>WASHINGTON DC</COUNTY>
    <PRIMARY_CARE_PROVIDER>HERSH, WILLIAM R.</PRIMARY_CARE_PROVIDER>
    <BIOLOGY_SAMPLE_OPT_OUT>N</BIOLOGY_SAMPLE_OPT_OUT>
    <GENETIC_OPT_OUT>N</GENETIC_OPT_OUT>
  </DATA_RECORD>
  - <DATA_RECORD>
    <MEDICAL_RECORD_NUMBER>11111111</MEDICAL_RECORD_NUMBER>
    <PATIENT_ID>Z1111111</PATIENT_ID>
    <CURRENT_AGE_IN_YEARS>15</CURRENT_AGE_IN_YEARS>
    <BIRTH_DATE>09/05/2003</BIRTH_DATE>
    <GENDER>FEMALE</GENDER>
    <ETHNICITY>NOT HISPANIC OR LATINO</ETHNICITY>
    <RACE>ASIAN</RACE>
    <ALIVE>Y</ALIVE>
    <DEATH_DATE>12/31/1999</DEATH_DATE>
    <LAST_NAME>HUQUE</LAST_NAME>
    <FIRST_NAME>ALAYNA</FIRST_NAME>
    <CITY>WASHINGTON</CITY>
    <ADDRESS_LINE1>1600 PENNSYLVANIA AVENUE</ADDRESS_LINE1>
    <ADDRESS_LINE2>N/A</ADDRESS_LINE2>
    <STATE>DC</STATE>
    <ZIP>20500</ZIP>
    <PHONE>2024561414</PHONE>
    <EMAIL>ALAHUQUE&#64GMAIL.COM</EMAIL>
    <COUNTY>WASHINGTON DC</COUNTY>
    <PRIMARY_CARE_PROVIDER>HERSH, WILLIAM R.</PRIMARY_CARE_PROVIDER>
    <BIOLOGY_SAMPLE_OPT_OUT>N</BIOLOGY_SAMPLE_OPT_OUT>
    <GENETIC_OPT_OUT>N</GENETIC_OPT_OUT>
  </DATA_RECORD>
</main>
```

# Python Program

1. Parse the OHSU XML source files
2. Store the EHR patient data into a list
3. Add all lists to a queue
4. Remove one list at a time from the queue
5. Index the list and assign the information to the OMOP field
6. Print a new XML file that follows the OMOP CDM standards

## Steps 1 and 2

```
1 def demographics(data_set):
2     """Parses the xml file for the demographics table of patients and
3     returns a list of lists containing the data for the OMOP CDM table"""
4     base_path = os.path.dirname(os.path.realpath(__file__))
5     xml_file = os.path.join(base_path, "C:\Users\hugue\Desktop\data\{}", format(data_set))
6     tree = et.parse(xml_file) # parses xml file
7     root = tree.getroot()
8
9     # creates a list for each DATA_RECORD
10    collected_data = []
11    for d in root.findall('DATA_RECORD'):
12        collected_data.append([d.find(x).text for x in ('OHSU_PATIENT_ID',
13            'BIRTH_DATE',
14            'GENDER',
15            'ETHNICITY',
16            'RACE',
17            'DEATH_DATE',
18            'LAST_NAME',
19            'FIRST_NAME',
20            'CITY',
21            'ADDRESS_LINE1',
22            'ADDRESS_LINE2',
23            'STATE',
24            'ZIP',
25            'COUNTY',
26            'PRIMARY_CARE_PROVIDER')])
27
28    return collected_data
```

## Steps 3 and 4

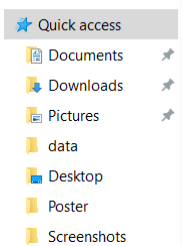
```
1 def o_demographics(data_set):
2     """Adds the list containing data records of demographics
3     to a queue and continuously removes a list and prints an OMOP PERSON
4     and LOCATION xml document until the queue is empty"""
5     demographics_queue = Queue()
6     for i in range(len(gphs.demographics(data_set))):
7         demographics_queue.put(demographics(data_set[i])) # adds lists to queue
8     while demographics_queue.isempty() != True:
9         data_record = demographics_queue.get() # removes list from queue
10        gph.print_demographics_PERSON_LOCATION(data_record) # prints OMOP xml document
```





## Step 5

```
1 def demographics_PERSON_elements(root, collected_data):
2     """Indexes the demographics list and maps the OMOP data fields to
3     the OMOP structure for the PERSON table"""
4     record = et.SubElement(root, 'DATA_RECORD')
5     et.SubElement(record, 'PERSON_ID').text = collected_data[0]
6     et.SubElement(record, 'gender_source_value').text = collected_data[1]
7     et.SubElement(record, 'year_of_birth').text = collected_data[10]
8     et.SubElement(record, 'month_of_birth').text = collected_data[11]
9     et.SubElement(record, 'day_of_birth').text = collected_data[12]
10    et.SubElement(record, 'birth_datetime').text = collected_data[13]
11    et.SubElement(record, 'death_datetime').text = collected_data[14]
12    et.SubElement(record, 'race_source_value').text = collected_data[15]
13    et.SubElement(record, 'ethnicity_source_value').text = collected_data[16]
14    et.SubElement(record, 'location_id').text = collected_data[17]
15    et.SubElement(record, 'provider_id').text = collected_data[18]
16    et.SubElement(record, 'care_site_id').text = collected_data[19]
17    et.SubElement(record, 'person_source_value').text = collected_data[20]
18    et.SubElement(record, 'gender_source_value').text = collected_data[21]
19    et.SubElement(record, 'gender_source_concept_id').text = collected_data[22]
20    et.SubElement(record, 'race_source_value').text = collected_data[23]
21    et.SubElement(record, 'race_source_concept_id').text = collected_data[24]
22    et.SubElement(record, 'ethnicity_source_value').text = collected_data[25]
23    et.SubElement(record, 'ethnicity_source_concept_id').text = collected_data[26]
24
25 def demographics_LOCATION_elements(root, collected_data):
26     """Indexes the demographics list and maps the OMOP data fields to
27     the OMOP structure for the LOCATION table"""
28     record = et.SubElement(root, 'DATA_RECORD')
29     et.SubElement(record, 'location_id').text = collected_data[17]
30    et.SubElement(record, 'address_1').text = collected_data[27]
31    et.SubElement(record, 'address_2').text = collected_data[28]
32    et.SubElement(record, 'city').text = collected_data[29]
33    et.SubElement(record, 'state').text = collected_data[30]
34    et.SubElement(record, 'zip').text = collected_data[31]
35    et.SubElement(record, 'country').text = collected_data[32]
36    et.SubElement(record, 'location_source_value').text = collected_data[33]
37    et.SubElement(record, 'location_source_concept_id').text = collected_data[34]
38    et.SubElement(record, 'latitude').text = collected_data[35]
39    et.SubElement(record, 'longitude').text = collected_data[36]
```

## Step 6

```
1 def print_demographics_PERSON_LOCATION(collected_data):
2     """Prints a list of xml documents for the OMOP PERSON and LOCATION
3     tables to a PERSON and LOCATION table that meets OMOP standards"""
4     root = et.Element('LOCATION')
5     gph.demographics_PERSON_elements(root, collected_data) # indexes list
6
7     tree = et.ElementTree(root)
8     tree.write("C:\Users\hugue\Desktop\data\PERSON_collected_data[0].xml")
9
10    root = et.Element('PERSON')
11    gph.demographics_LOCATION_elements(root, collected_data) # indexes list
12
13    tree = et.ElementTree(root)
14    tree.write("C:\Users\hugue\Desktop\data\LOCATION_collected_data[0].xml")
```



	LOCATION_Z0000000	Da	Siz
	LOCATION_Z1111111	Da	Siz
	PERSON_Z0000000	Da	Siz
	PERSON_Z1111111	Da	Siz

# Example OMOP EHR XML Document PERSON and LOCATION

```
<?xml version="1.0"?>
- <PERSON>
  - <DATA_RECORD>
    <person_id>Z0000000</person_id>
    <gender_concept_id/>
    <year_of_birth>1999</year_of_birth>
    <month_of_birth>04</month_of_birth>
    <day_of_birth>21</day_of_birth>
    <birth_datetime>04/21/1999</birth_datetime>
    <death_datetime>12/31/1999</death_datetime>
    <race_concept_id/>
    <ethnicity_concept_id/>
    <location_id>1</location_id>
    <provider_id/>
    <care_site_id/>
    <person_source_value>ALYSSA HUQUE</person_source_value>
    <gender_source_value>FEMALE</gender_source_value>
    <gender_source_concept_id/>
    <race_source_value>ASIAN</race_source_value>
    <race_source_concept_id/>
    <ethnicity_source_value>NOT HISPANIC OR LATINO</ethnicity_source_value>
    <ethnicity_source_concept_id/>
  </DATA_RECORD>
</PERSON>
```

```
<?xml version="1.0"?>
- <LOCATION>
  - <DATA_RECORD>
    <location_id>2</location_id>
    <address_1>1600 PENNSYLVANIA AVENUE</address_1>
    <address_2>N/A</address_2>
    <city>WASHINGTON</city>
    <state>DC</state>
    <zip>20500</zip>
    <county>WASHINGTON DC</county>
    <country/>
    <location_source_value/>
    <latitude/>
    <longitude/>
  </DATA_RECORD>
</LOCATION>
```



# The Two Programs

- Similar approaches in both programs
  - Was able to map and standardize to OMOP CDM
- Programs can be used in tandem
  - Print OMOP XML documents
  - Move into database (such as SQL)
  - Implement standardized codes and relationality of tables



# Future Directions

- Refine OHSU to OMOP mappings
- Implement standardized codes used in the OMOP CDM `_concept_id` fields
- Improve the efficiency of our programs to handle a few million data records



Funded by Grant 2T15LM007088, Research Training in Biomedical  
Informatics and Data Science at Oregon Health and Science University  
National Library of Medicine

**Special thanks to Dr. William Hersh, Dr. Aaron Cohen,  
Dr. Steven Bedrick, and Dr. Steven Chamberlin**



# Works Cited

- Belenkaya, R., Blacketer, C., Hripcsak, G., Natajan, K., O'Hara, D., Reich, D., Roa, G., Torok, D., Zandt, M., Velez, M., Voss, E.. *OMOP Common Data Model and Standardized Vocabularies*. <https://www.ohdsi.org/wp-content/uploads/2017/10/OHDSI-Vocabulary-CDM-Tutorial-2017-2017.10.11.pdf>. PowerPoint Presentation.
- Ohdsi. "OHDI/CommonDataModel." *GitHub*, [github.com/OHDSI/CommonDataModel/wiki/Standardized-Clinical-Data-Tables](https://github.com/OHDSI/CommonDataModel/wiki/Standardized-Clinical-Data-Tables).
- "OMOP Common Data Model." *OHDSI*, [www.ohdsi.org/data-standardization/the-common-data-model](http://www.ohdsi.org/data-standardization/the-common-data-model). Accessed 26 Jun. 2019.
- "Rodriguez, Alison. "HITECH Act Resulted in Significant Gains in EHR Adoption in Hospitals." *AJMC*, 15 Aug. 2017, [www.ajmc.com/newsroom/hitech-act-resulted-in-significant-gains-in-ehr-adoption-in-hospitals](http://www.ajmc.com/newsroom/hitech-act-resulted-in-significant-gains-in-ehr-adoption-in-hospitals).