# OMOP
# Common Data Model (CDM)
# & Extract-Transform-Load (ETL)
# Tutorial

Rimma Belenkaya (Memorial Sloan Kettering)
Karthik Natarajan (Columbia University)
Mark Velez (Columbia University)
Erica Voss (Janssen R&D Epidemiology Analytics)

24 September 2016

**Please copy the contents of the
USB drive to your hard disk now.
You will need ~45GB free disk space available.**

| Rimma Belenkaya<br>Memorial Sloan Kettering | Karthik Natarajan<br>Columbia University |
| --- | --- |
| | |
| Mark Velez<br>Columbia University | Erica Voss<br>Janssen R&D Epidemiology Analytics |
| | |

# Teaching Assistants

| Anthony Sena<br>Janssen R&D Epidemiology Analytics | Jungmi Han<br>Columbia University |
|:---:|:---:|
|  |  |

# Ground Rules

- We are recording today's session, so presenters should repeat questions.

- We may table source specific questions.

- The Virtual Machine (VM) distributed today on USB, please return.

- If we cannot get the VM working on your machine let's try to buddy you up.  Do not worry the presentation will still walk you through the content.

- This course will not focus on the Vocabulary, however the Vocabulary is critical to the Common Data Model and the ETL process.

# Agenda

| Time | Type | Section |
|------|------|---------|
| 8:00AM-8:15AM | | Introductions |
| 8:15AM-9:15AM | *Foundational* | What is OMOP/OHDSI? OMOP Common Data Model (CDM) – Why and How |
| 9:15AM-10:00AM | | How to retrieve data from OMOP CDM |
| 10:00AM-10:15AM | | Break |
| 10:15AM-10:45AM | *Implementation* | Setup and Performing of an Extract Transform and Load process into the CDM |
| 10:45AM-11:30AM | | Using WhiteRabbit and Rabbit-In-A-Hat to Build an ETL |
| 11:30AM-11:45AM | *Evaluation* | Testing and Quality Assurance |
| 11:45AM-12:00PM | | Wrap up |

# **Foundational**

What is OMOP/OHDSI?
OMOP Common Data Model
(CDM) – Why and How

# Introduction of OMOP/OHDSI

**OHDSI: Observational Health Data Sciences and Informatics** is a research collaborative coordinated through Columbia University

**Who?**

–Multiple stakeholders: academia, government, industry

–Multiple disciplines: statistics, epidemiology, informatics, clinical sciences

**Why?** To generate evidence about all aspects of healthcare

**Where?** Multiple geographies: US, Europe, Asia-Pacific, 20 countries. OHDSI collaborators access a network of 600 mln patients

**How?** By developing analytical methods and tools based on the data standardized to OMOP Common Data Model (CDM) and vocabulary

# OMOP Common Data Model (CDM) What is it and why have one?

**What?**

- A standardized way to represent data structure (CDM) and content (vocabulary)
- One model to accommodate data coming from disparate data sources
    - –administrative claims, electronic health records
    - –EHRs from both inpatient and outpatient settings
    - –registries and longitudinal surveys
    - –data sources both within and outside of US

**Why?**

- Enable standardization of structure and content to support a systematic and reproducible process to efficiently generate evidence
- Support collaborative research both within and outside of US

# OMOP CDM v5.0.1

# OMOP CDM Design Principles

- Relational design but platform independent
  - Integrated with Controlled Vocabulary
  - Domain (subject area) based
  - Patient centric
  - Uniformly integrates data from heterogeneous data sources: EMR, claims, registries
- Built for analytical purposes, extended/developed based on analytic use cases
- Extendable, both vocabulary (new vocabs, local concepts) and CDM (Observation)
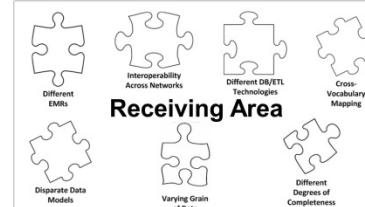
# NYC-CDRN Experience

# OMOP CDM v5.0.1

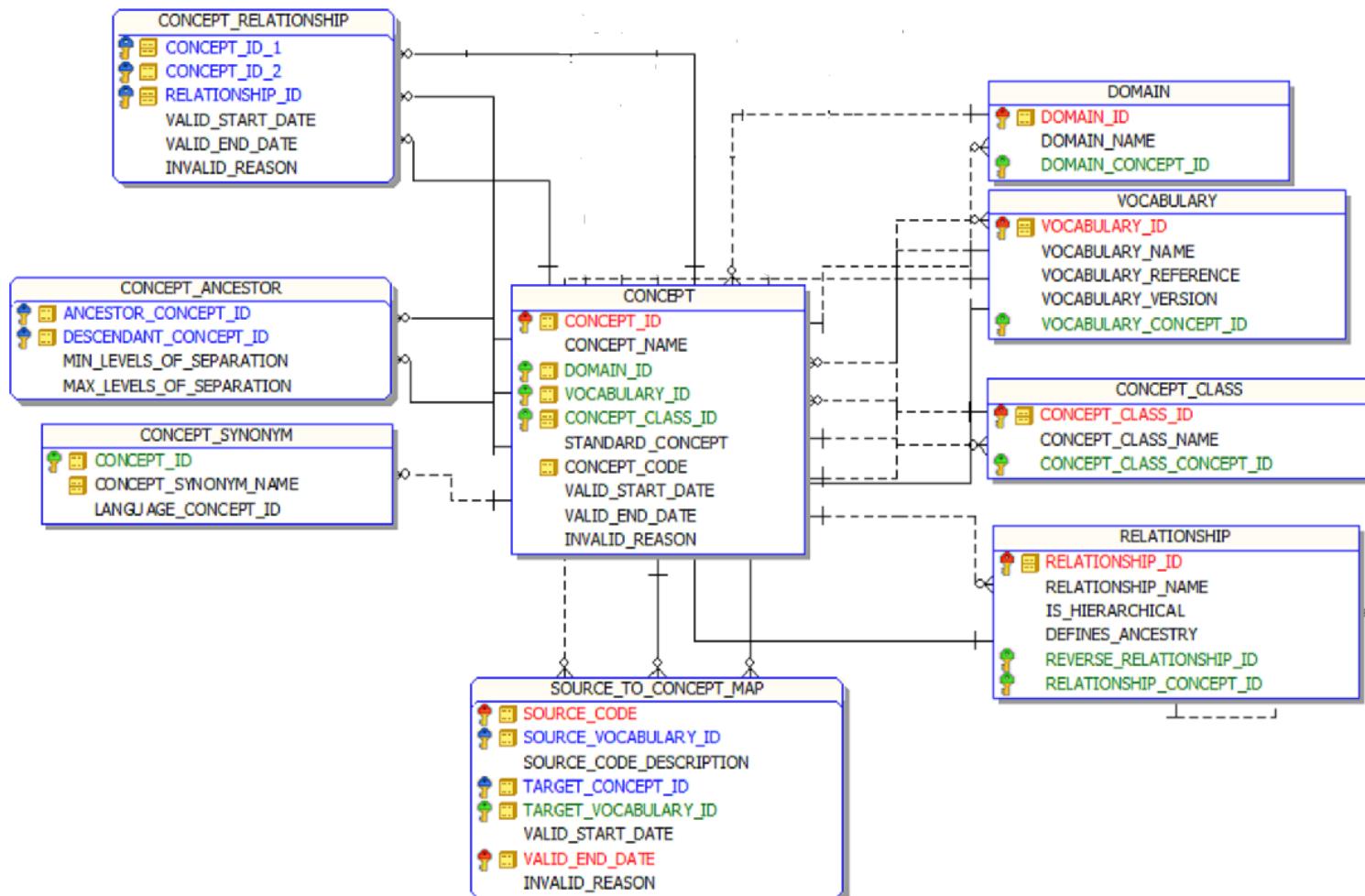# OMOP Common Vocabulary Model

**What it is**

- Standardized structure to house existing vocabularies used in the public domain
- Compiled standards from disparate public and private sources and some OMOP-grown concepts
- Built on the shoulders of National Library of Medicine's Unified Medical Language System (UMLS)

**What it's not**

- Static dataset – the vocabulary updates regularly to keep up with the continual evolution of the sources
- Finished product – vocabulary maintenance and improvement is ongoing activity that requires community participation and support
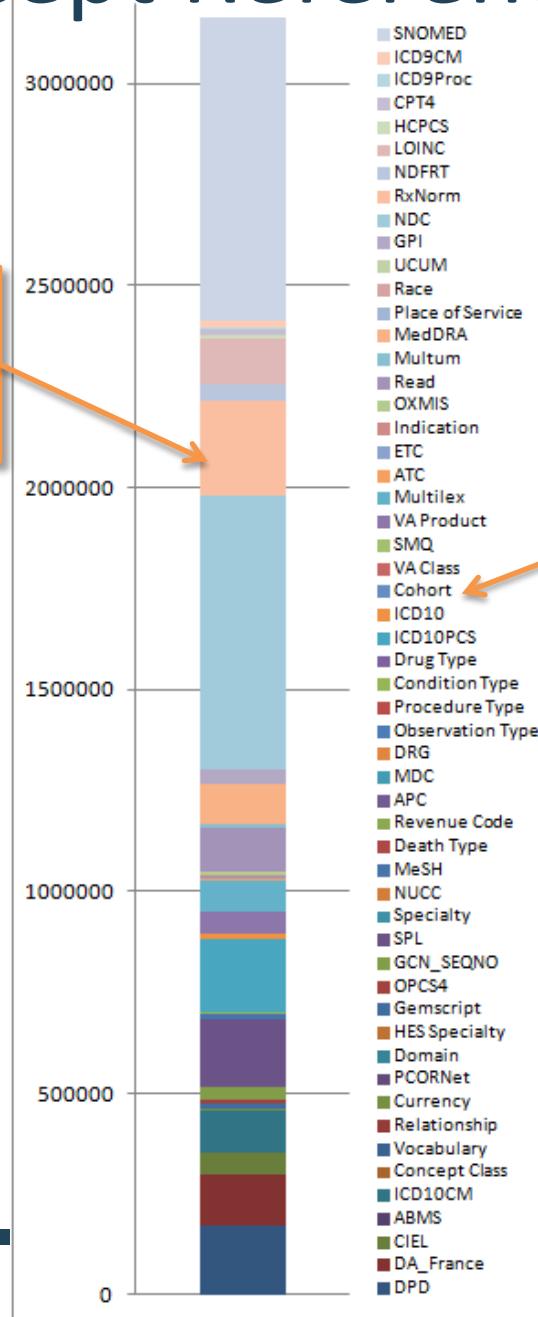
# OMOP Common Vocabulary Model

# Single Concept Reference Table



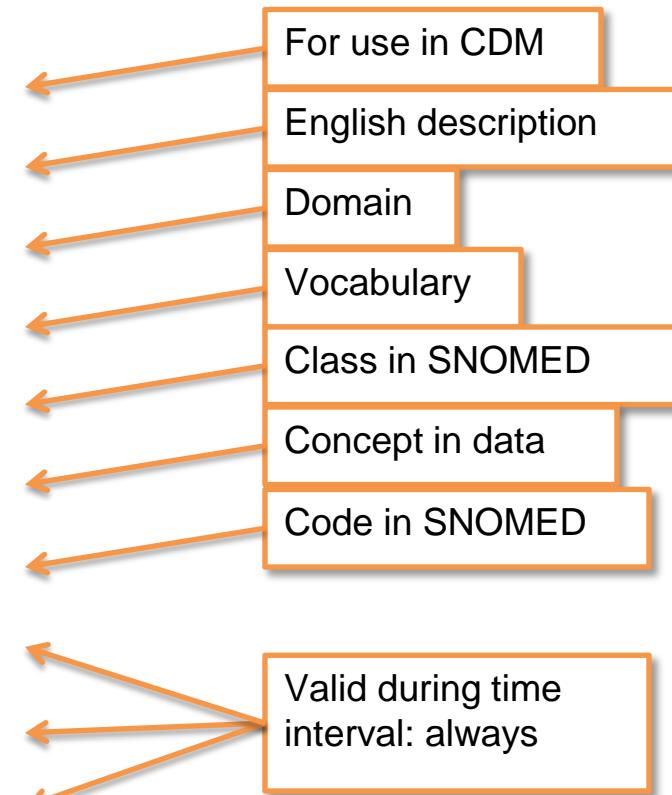All vocabularies stacked up in one table

Vocabulary ID

# What's in a Concept

| | |
|---|---|
| CONCEPT_ID | 313217 |
| CONCEPT_NAME | Atrial fibrillation |
| DOMAIN_ID | Condition |
| VOCABULARY_ID | SNOMED |
| CONCEPT_CLASS_ID | Clinical Finding |
| STANDARD_CONCEPT | S |
| CONCEPT_CODE | 49436004 |
| VALID_START_DATE | 01-Jan-70 |
| VALID_END_DATE | 31-Dec-99 |
| INVALID_REASON | |

For use in CDM

English description

Domain

Vocabulary

Class in SNOMED

Concept in data

Code in SNOMED

Valid during time interval: always

# OMOP Vocabulary Model Design Principles

- Uniform structure
  - All concepts are in one table
  - All concept relationships are in one table, including mappings from source to standard vocabularies
- Formalized integration with Common Data Model via concept domain
  - Direction of ETL is informed by concept domain
- Relationships are bi-directional
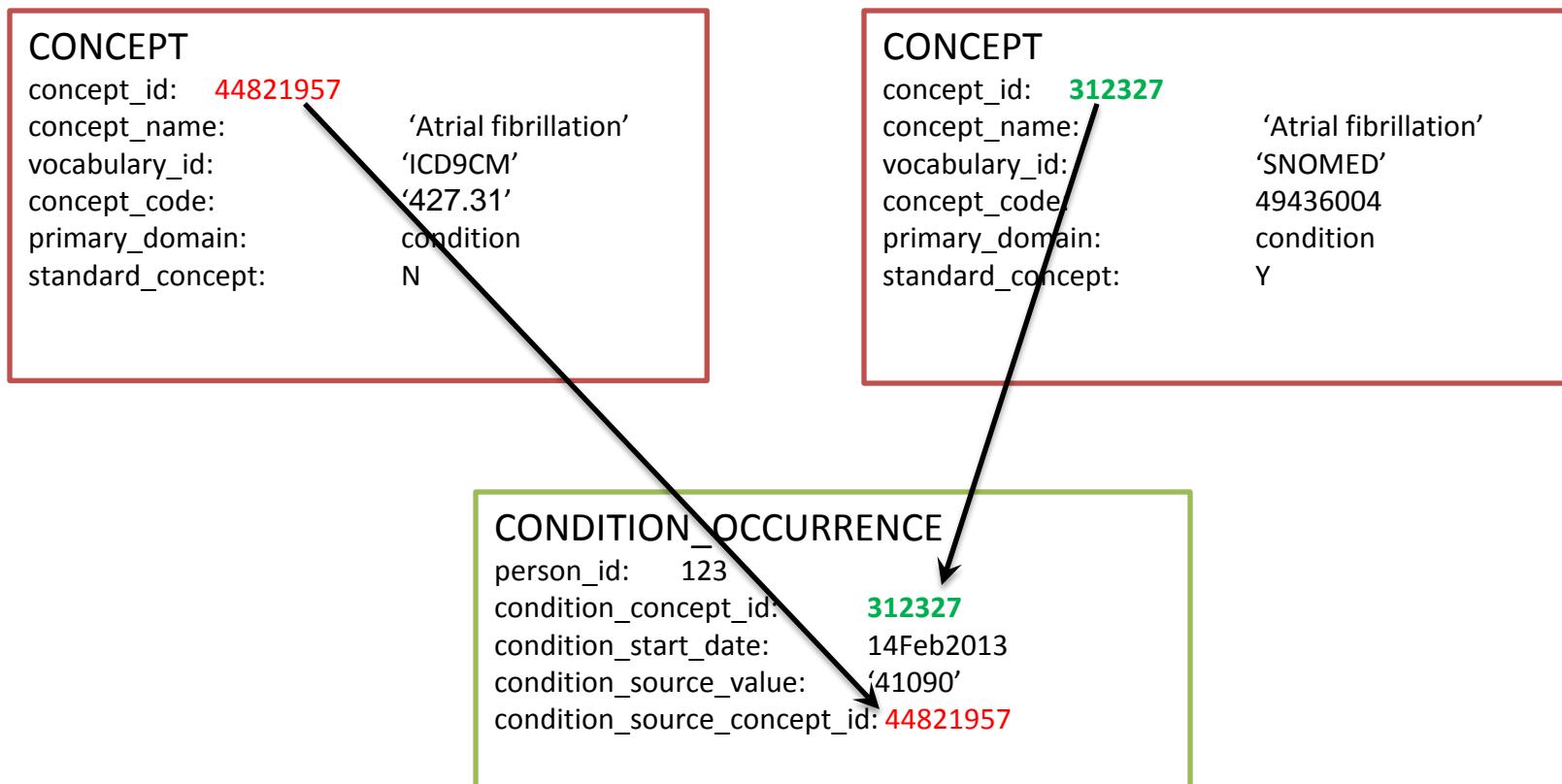- Hierarchical relationships have additional representation in the model to support efficient data retrieval
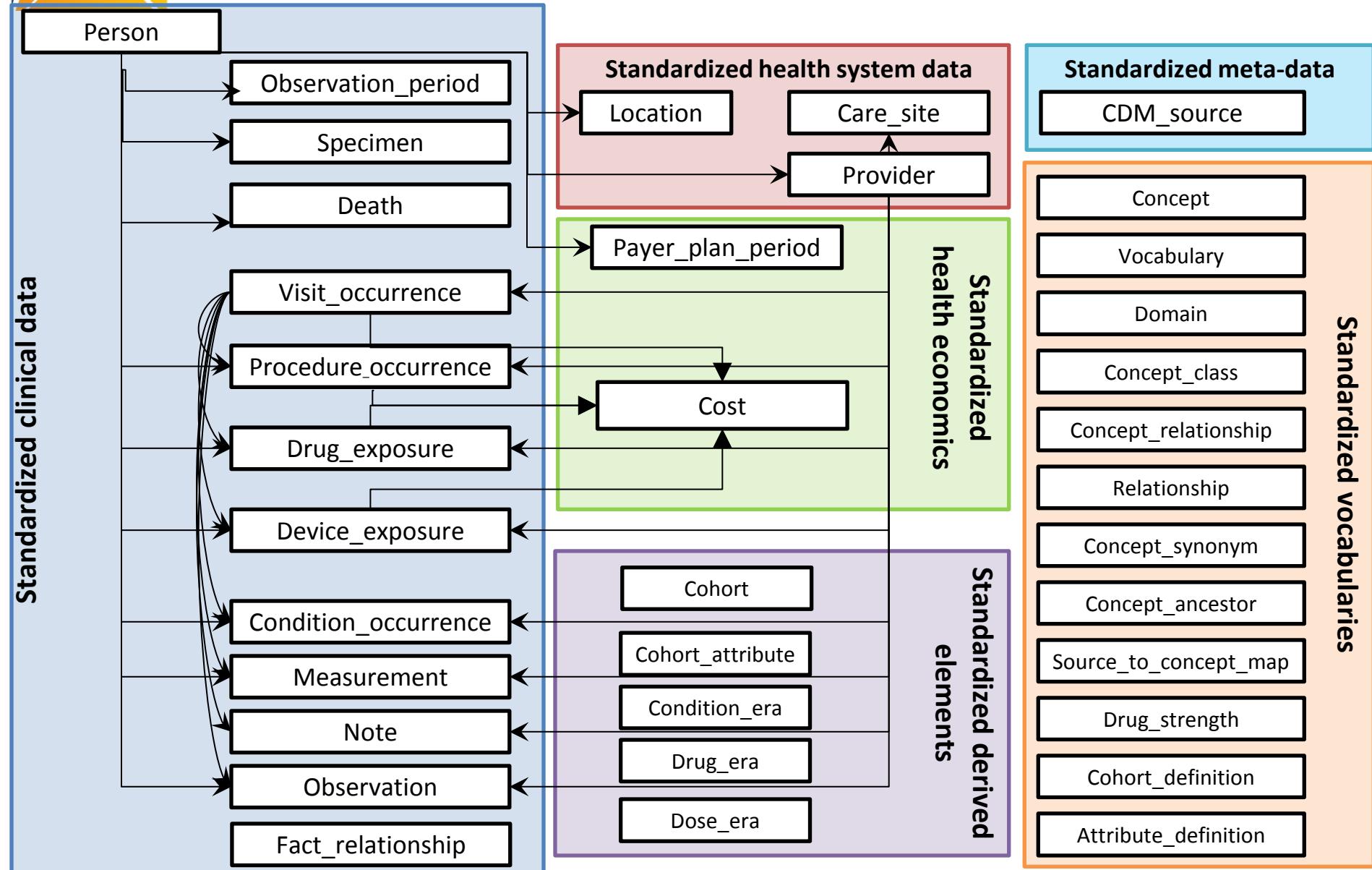
# OMOP CDM Standard Domain Features

| Feature | Description and purpose | Field name convention | Example |
|---|---|---|---|
| Patient centric | Every domain table has **patient identifier**. Patient data can be retrieved independently from other domains. | **person_id** | person_id  123 |
| Unique domain identifier | Every domain table has a unique primary key to identify domain **entities** | <entity>**_id** | condition_occurrence_id 470985 |
| Standard concept from a respective vocabulary domain | Integration with the vocabulary. Foreign key into the Standard Vocabulary for **Standard Concept** | <entity>**_concept_id** | condition_concept_id 313217 (SNOMED "Atrial Fibrillation") |
| Source concept from a respective vocabulary domain | Provenance. Foreign key into the Standard Vocabulary for **Source Concept** | <entity>**_source_concept_id** | condition_source_concept_id 44821957 (ICD9CM "Atrial Fibrillation") |
| Source value | Provenance. Verbatim information from the source data, **not to be used** by any standard analytics | <entity>**_source_value** | condition_source_value 427.31 (ICD9CM "Atrial Fibrillation") |
| Source type | Provenance. Foreign key into the Vocabulary for the **origin of the** | <entity>**_type_concept_id** | condition_type_concept_id 38000199 ("Inpatient header – primary") |

# Integration of CDM and Vocabulary

**CONCEPT**
concept_id:          44821957
concept_name:                  'Atrial fibrillation'
vocabulary_id:        'ICD9CM'
concept_code:         '427.31'
primary_domain:       condition
standard_concept:     N

**CONCEPT**
concept_id:          312327
concept_name:                  'Atrial fibrillation'
vocabulary_id:        'SNOMED'
concept_code:         49436004
primary_domain:       condition
standard_concept:     Y

**CONDITION_OCCURRENCE**
person_id:       123
condition_concept_id:          312327
condition_start_date:          14Feb2013
condition_source_value:        '41090'
condition_source_concept_id:   44821957

# OMOP CDM v5.0.1

# PERSON

| person | |
|---|---|
| 🔑 | person_id |
| | gender_concept_id |
| | year_of_birth |
| | month_of_birth |
| | day_of_birth |
| | time_of_birth |
| | race_concept_id |
| | ethnicity_concept_id |
| | location_id |
| | provider_id |
| | care_site_id |
| | person_source_value |
| | gender_source_value |
| | gender_source_concept_id |
| | race_source_value |
| | race_source_concept_id |
| | ethnicity_source_value |
| | ethnicity_source_concept_id |

- Need to create one unique record per person (not multiple rows per move)

- Vocabulary for gender, race, ethnicity: HL7 administrative

- No history of location/demographics: need to select latest available

- Location peculiarity: foreign key to the LOCATION table that contains one record per each unique location

- Year of birth required…day/month optional

# LOCATION

| location | |
|---|---|
| 🔑 | location_id |
| | address_1 |
| | address_2 |
| | city |
| | state |
| | zip |
| | county |
| | location_source_value |

- Contains one record per each unique location

- Location is highly variable across sources, of limited use thus far

# OBSERVATION_PERIOD

**observation_period**

| | |
|---|---|
| 🔑 | observation_period_id |
| | person_id |
| | observation_period_start_date |
| | observation_period_end_date |
| | period_type_concept_id |

- Spans of time where data source has capture of data

- Required to run analytical methods

- One person may have multiple periods if there is interruption in data capture

- Challenge: determine observation periods based on the source data

# DEATH

| death | |
|---|---|
| 🔑 | person_id |
| | death_date |
| | death_type_concept_id |
| | cause_concept_id |
| | cause_source_value |
| | cause_source_concept_id |

- Can have death without cause

- Can only have 1 death per person

# VISIT_OCCURRENCE

**visit_occurrence**

- visit_occurrence_id
- person_id
- visit_concept_id
- visit_start_date
- visit_start_time
- visit_end_date
- visit_end_time
- visit_type_concept_id
- provider_id
- care_site_id
- visit_source_value
- visit_source_concept_id

- Visits <> 'Encounters':
  - claims often need to be consolidated to minimize double-counting
  - inpatient transitions are not covered

- Visit Types
  - Inpatient
  - Emergency room
  - Inpatient/Emergency - new
  - Outpatient
  - Long-term care

- Vocabulary: OMOP

- Other attributes: time of visit start/end, provider, admitting source, discharge disposition

# PROCEDURE_OCCURRENCE

| procedure_occurrence |
| --- |
| procedure_occurrence_id |
| person_id |
| procedure_concept_id |
| procedure_date |
| procedure_type_concept_id |
| modifier_concept_id |
| quantity |
| provider_id |
| visit_occurrence_id |
| procedure_source_value |
| procedure_source_concept_id |
| qualifier_source_value |

- Vocabularies: CPT-4, HCPCS, ICD-9 Procedures, ICD-10 Procedures, LOINC, SNOMED

- Procedures have the least standardized vocabularies that causes some redundancy

# CONDITION_OCCURRENCE

**condition_occurrence**

| | |
|---|---|
| 🔑 | condition_occurrence_id |
| | person_id |
| | condition_concept_id |
| | condition_start_date |
| | condition_end_date |
| | condition_type_concept_id |
| | stop_reason |
| | provider_id |
| | visit_occurrence_id |
| | condition_source_value |
| | condition_source_concept_id |

- Vocabulary: SNOMED -> classification

- Data sources:
  - Billing diagnosis (inpatient, outpatient)
  - Problem list

- Individual records <> distinct episodes

# DRUG_EXPOSURE

**drug_exposure**

- 🔑 drug_exposure_id
- person_id
- drug_concept_id
- drug_exposure_start_date
- drug_exposure_end_date
- drug_type_concept_id
- stop_reason
- refills
- quantity
- days_supply
- sig
- route_concept_id
- effective_drug_dose
- dose_unit_concept_id
- lot_number
- provider_id
- visit_occurrence_id
- drug_source_value
- drug_source_concept_id
- route_source_value
- dose_unit_source_value

- Vocabulary: RxNorm-> classifications by drug class and indication

- Data sources:
  - Pharmacy dispensing
  - Prescriptions written
  - Medication history

- Source fields may vary, but so inference of drug exposure end may vary

# DEVICE_EXPOSURE

**device exposure**

| | |
|---|---|
| 🔑 | device_exposure_id |
| | person_id |
| | device_concept_id |
| | device_exposure_start_date |
| | device_exposure_end_date |
| | device_type_concept_id |
| | unique_device_id |
| | quantity |
| | provider_id |
| | visit_occurrence_id |
| | device_source_value |
| | device_source_concept_id |

- OMOP CDM is the only data model supporting devices
- Accommodates FDA unique device identifiers (UDI) even though most data sources don't have them yet

# MEASUREMENT

**measurement**

- measurement_id
- person_id
- measurement_concept_id
- measurement_date
- measurement_time
- measurement_type_concept_id
- operator_concept_id
- value_as_number
- value_as_concept_id
- unit_concept_id
- range_low
- range_high
- provider_id
- visit_occurrence_id
- measurement_source_value
- measurement_source_concept...
- unit_source_value
- value_source_value

- EAV design

- Vocabulary: LOINC, SNOMED

- Data sources: structured, quantitative measures, such as laboratory tests

- Measures have associated units
  - Measurement units vocabulary: UCUM

- No free format for measurement results

# OBSERVATION

| observation | |
|---|---|
| 🔑 | observation_id |
| | person_id |
| | observation_concept_id |
| | observation_date |
| | observation_time |
| | observation_type_concept_id |
| | value_as_number |
| | value_as_string |
| | value_as_concept_id |
| | qualifier_concept_id |
| | unit_concept_id |
| | provider_id |
| | visit_occurrence_id |
| | observation_source_value |
| | observation_source_concept_id |
| | unit_source_value |
| | qualifier_source_value |

- Catch-all EAV design to capture all other data:
  - observation: 'question'
  - value: 'answer'
    - Can be numeric, concept, or string (e.g. free text)

- Instrument for CDM extension, playpen

- Not all 'questions' are standardized, source value can accommodate 'custom' observations (particularly pertinent in registries)

# SPECIMEN

| specimen |
| --- |
| 🔑 specimen_id |
| person_id |
| specimen_concept_id |
| specimen_type_concept_id |
| specimen_date |
| specimen_time |
| quantity |
| unit_concept_id |
| anatomic_site_concept_id |
| disease_status_concept_id |
| specimen_source_id |
| specimen_source_value |
| unit_source_value |
| anatomic_site_source_value |
| disease_status_source_value |

- To capture of biomarker / tissue bank

# NOTE

| note |  |
|---|---|
| 🔑 | note_id |
|  | person_id |
|  | note_date |
|  | note_time |
|  | note_type_concept_id |
|  | note_text |
|  | provider_id |
|  | note_source_value |
|  | visit_occurrence_id |

- To capture unstructured free text

- Coming soon in CDM 5.x: NLP and LOINC Clinical Document Ontology (CDO) annotations

# Health Economics

**payer_plan_period**

| | |
|---|---|
| 🔑 | payer_plan_period_id |
| | person_id |
| | payer_plan_period_start_date |
| | payer_plan_period_end_date |
| | payer_source_value |
| | plan_source_value |
| | family_source_value |

**cost**

| | |
|---|---|
| 🔑 | cost_id |
| | cost_event_id |
| | cost_domain_id |
| | cost_type_concept_id |
| | currency_concept_id |
| | total_charge |
| | total_cost |
| | total_paid |
| | paid_by_payer |
| | paid_by_patient |
| | paid_patient_copay |
| | paid_patient_coinsurance |
| | paid_patient_deductible |
| | paid_by_primary |
| | paid_ingredient_cost |
| | paid_dispensing_fee |
| | payer_plan_period_id |
| | amount_allowed |
| | revenue_code_concept_id |
| | revenue_code_source_value |

- All costs consolidated into one table COST table

- Costs tied to respective observation records

- Domain is determined by cost_domain_id (e.g. visit, condition, etc.)

# OMOP CDM Service Tables

- CDM_SOURCE
  - Provenance, integration, metadata
  - Future extension to individual domains

- FACT_RELATIONSHIP
  - Linkage between related observations
  - Example: systolic and diastolic blood pressure

# Motivation for Standardized Derived Elements

- Derived elements intended to supplement- not replace- raw data
  - If derived assumptions don't meet a specific use case, don't use them

- Promotes transparency and consistency in research by having standard processes applies across analyses

- Increased efficiency by processing key data elements once at ETL-time, rather than requiring each analysis to figure it out at each analysis run-time

- Key standardized elements available in OMOP CDMv5:
  - Cohort – standardize definition and syntax for defining populations that meet inclusion criteria
  - Drug era – standardize inference of length of exposure to product for all active ingredients
  - Dose era – standardize estimation of daily dose for periods of exposure to all drug products
  - Condition era – standardize aggregation of episodes of care, delineating between acute vs. chronic conditions

# Cohort Management



1. **COHORT** table contains records of subjects that satisfy a given set of criteria for a duration of time.

2. The definition of the cohort is contained within the **COHORT_DEFINITIO**N table. It provides a standardized structure for maintaining the rules governing the inclusion of a subject into a cohort, and can store programming code to instantiate the cohort within the OMOP CDM.

3. **COHORT_ATTRIBUTE** table contains attributes associated with each subject within a cohort, as defined by a given set of criteria for a duration of time.

4. The definition of the Cohort Attribute is contained in the **ATTRIBUTE_DEFINITION** table.

# DRUG_ERA

| drug_era | |
|---|---|
| 🔑 | drug_era_id |
| | person_id |
| | drug_concept_id |
| | drug_era_start_date |
| | drug_era_end_date |
| | drug_exposure_count |
| | gap_days |

- Standardized inference of length of exposure to product for all active ingredients

- Derived from records in DRUG_EXPOSURE under certain rules to produce continuous Drug Eras

Illustrating inferences needed within longitudinal pharmacy claims data for one patient

# What makes OMOP CDM unique

- Specialized CDM - reflective of clinical domain, granular, well structured

- Vocabulary - uniformly structured and well curated

- Information Model - formalized connection between data model and conceptual model (Vocabulary)

- Specialized yet Extendable – new attributes and concepts can be added

- Supportive Community of developers and researchers

- Development driven by analytic use cases

# Foundational

## How to retrieve data from OMOP CDM

# OHDSI in a Box

VirtualBox

**PostgreSQL**

cdm    webapi

pgAdmin

**docker**

**Broadsea**

**WebTools**

| Atlas | Penelope | Calypso |

WebAPI

Tomcat

**Methods Library**

| | OHDSI R packages | Studio |

synpuf_100k

WhiteRabbit

RabbitInAHat

# OHDSI in a Box – Setup

1. Open ▣ VM VirtualBox Manager

2. Click on ✶
   New

## Name and operating system

Please choose a descriptive name for the new virtual machine and select the type of operating system you intend to install on it. The name you choose will be used throughout VirtualBox to identify this machine.

Name: OHDSI-1percent

Type: Linux

Version: Ubuntu (64-bit)

## Memory size

Select the amount of memory (RAM) in megabytes to be allocated to the virtual machine.

The recommended memory size is **1024** MB.

2048 MB

4 MB                    8192 MB

○ Do not add a virtual hard disk

○ Create a virtual hard disk now

● Use an existing virtual hard disk file

ohdsi - 1k - Final.vdi (Normal, 30.00 GB)

Oracle VM VirtualBox Manager

File   Machine   Help

New   Settings   Discard   Start          Details   Snapshots

ohdsi-1p
⏻ Powered Off

**General**
Name:              ohdsi-1p
Operating System:  Ubuntu (64-bit)

**System**
Base Memory:   2048 MB
Boot Order:    Floppy, Optical, Hard Disk
Acceleration:  VT-x/AMD-V, Nested Paging, KVM Paravirtualization

**Preview**

ohdsi-1p

**Display**
Video Memory:           16 MB
Remote Desktop Server:  Disabled
Video Capture:          Disabled

**Storage**
Controller: IDE
  IDE Secondary Master:  [Optical Drive] Empty
Controller: SATA
  SATA Port 0:           ohdsi_1percent.vdi (Normal, 50.29 GB)

**Audio**

# OHDSI in a Box – Start Up

# OHDSI in a Box – International Keyboards

# OHDSI in a Box – Adjust Resolution

# OHDSI in a Box – Clipboard

# OHDSI in a Box – Timeout

# OHDSI in a Box – Ready

# CDM Database – pgAdmin III New Server

# CDM Database – Connect

# CDM Database – Open SQL Sheet

# CDM Database – Ready

# Data Used for Demonstration

- Medicare Claims Synthetic Public Use Files (SynPUFs)
  - synthetic US Medicare insurance claims database
  - Medicare is a government based insurance program for primarily 65 and older but also individuals with disabilities
  - SynPUF not for research but rather demonstration/development purposes
  - Has been converted to the Common Data Model

**CMS.gov**
Centers for Medicare & Medicaid Services

# Data Used for Demonstration

- Five types of data:

| | DE-SynPUF | Unit of record | Number of Records 2008 | Number of Records 2009 | Number of Records 2010 |
|---|---|---|---|---|---|
| 1 | *Beneficiary Summary* | Benefi-ciary | 2,326,856 | 2,291,320 | 2,255,098 |
| 2 | *Inpatient Claims* | claim | 547,800 | 504,941 | 280,081 |
| 3 | *Outpatient Claims* | claim | 5,673,808 | 6,519,340 | 3,633,839 |
| 4 | *Carrier Claims* | claim | 34,276,324 | 37,304,993 | 23,282,135 |
| 5 | *Prescription Drug Events (PDE)* | event | 39,927,827 | 43,379,293 | 27,778,849 |

# SynPUF High Level Diagram

# Mapping SynPUF to CDM

# Some Example Questions

**Ex 1**

New Users of Warfarin

**Ex 2**

New Users of Warfarin
who are >=65?

**Ex 3**

New Users of Warfarin
with prior Atrial Fibrillation?

# New Users of Warfarin

- Warfarin is a blood thinner that is used to treat/prevent blood clots.

  – Where do you find drug data in the CDM?

  – What codes do I use to define drugs?

  – What does "New User" mean?

# Where are Drug Exposures in the CDM?



**Standardized clinical data**

- Person
- Observation_period
- Specimen
- Death
- Visit_occurrence
- Procedure_occurrence
- Drug_exposure
- Device_exposure
- Condition_occurrence
- Measurement
- Note
- Observation
- Fact_relationship

**Standardized health system data**

**Standardized meta-data**

**Standardized economics**

- Cost

**Standardized derived elements**

- Cohort
- Cohort_attribute
- Condition_era
- Drug_era
- Dose_era

**Standardized vocabularies**

- Concept_class
- Concept_relationship
- Relationship
- Concept_synonym
- Concept_ancestor
- Source_to_concept_map
- Drug_strength
- Cohort_definition
- Attribute_definition

captures records about the utilization of a drug when ingested or otherwise introduced into the body

# How do I define Warfarin?

**Ex 1**

- When raw data is transformed into the CDM raw source codes are transformed into standard OMOP Vocabulary concepts

- In the CDM, we no longer care what source concepts existed in the raw data, we just need to use concept identifiers

- We can use the OMOP Vocabulary to identify all concepts that contain the ingredient warfarin

# How do I define Warfarin?

**SQL**

- Writing SQL Statement

ATLAS

- OHDSI Tool ATLAS

# How do I define new users of a drug?

- someone who has recently started taking the drug, typically with a 6 or 12 month wash out



| 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 |

# How do I define new users of a drug?

- someone who has recently started taking the drug, typically with a 6 or 12 month wash out

index
drug

time in
database

6 months

**Ex 1**

# What is Needed in the CDM?

- **OMOP Vocabulary**
  to find the concepts

- **DRUG_EXPOSURE**
  to find individuals with exposure

- **OBSERVATION_PERIOD**
  to know people's time within the database

# New Users of Warfarin

```sql
/*****************************************************************************
 *     (Exercise 1) Warfarin New Users
 *****************************************************************************/

WITH CTE_DRUG_INDEX AS (
    SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
    FROM DRUG_EXPOSURE de
    WHERE de.DRUG_CONCEPT_ID IN (
        SELECT DESCENDANT_CONCEPT_ID
        FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
    )
    GROUP BY de.PERSON_ID
)
SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
    (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
FROM CTE_DRUG_INDEX i
    JOIN OBSERVATION_PERIOD op
        ON op.PERSON_ID = i.PERSON_ID
        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
ORDER BY i.PERSON_ID
```

# Step 1: Get the codes you need

```sql
/*************************************************************************
 *    (Exercise 1) Warfarin New Users
 *************************************************************************/

WITH CTE_DRUG_INDEX AS (
    SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
    FROM DRUG_EXPOSURE de
    WHERE de.DRUG_CONCEPT_ID IN (
        SELECT DESCENDANT_CONCEPT_ID
        FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
    )
    GROUP BY de.PERSON_ID
)
SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
    (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
FROM CTE_DRUG_INDEX i
    JOIN OBSERVATION_PERIOD op
        ON op.PERSON_ID = i.PERSON_ID
        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
ORDER BY i.PERSON_ID
```

# Step 2:  Find Drug Exposures

```
/********************************************************************
 *    (Exercise 1) Warfarin New Users
 ********************************************************************/

WITH CTE_DRUG_INDEX AS (
    SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
    FROM DRUG_EXPOSURE de
    WHERE de.DRUG_CONCEPT_ID IN (
        SELECT DESCENDANT_CONCEPT_ID
        FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
    )
    GROUP BY de.PERSON_ID
)
SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
    (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
FROM CTE_DRUG_INDEX i
    JOIN OBSERVATION_PERIOD op
        ON op.PERSON_ID = i.PERSON_ID
        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
ORDER BY i.PERSON_ID
```

# Step 3: Find New Users

```sql
/***********************************************************************
 *    (Exercise 1) Warfarin New Users
 **********************************************************************/

WITH CTE_DRUG_INDEX AS (
    SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
    FROM DRUG_EXPOSURE de
    WHERE de.DRUG_CONCEPT_ID IN (
        SELECT DESCENDANT_CONCEPT_ID
        FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
    )
    GROUP BY de.PERSON_ID
)
SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
    (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
FROM CTE_DRUG_INDEX i
    JOIN OBSERVATION_PERIOD op
        ON op.PERSON_ID = i.PERSON_ID
        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
ORDER BY i.PERSON_ID
```

# New Users of Warfarin

```
/****************************************************************************
 *    (Exercise 1) Warfarin New Users
 ****************************************************************************/

WITH CTE_DRUG_INDEX AS (
    SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
    FROM DRUG_EXPOSURE de
    WHERE de.DRUG_CONCEPT_ID IN (
        SELECT DESCENDANT_CONCEPT_ID
        FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
    )
    GROUP BY de.PERSON_ID
)
SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
    (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
FROM CTE_DRUG_INDEX i
    JOIN OBSERVATION_PERIOD op
        ON op.PERSON_ID = i.PERSON_ID
        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
ORDER BY i.PERSON_ID
```

# How do I define new users of warfarin who are >=65?

- someone who has recently started taking the drug, typically with a 6 or 12 month wash out

index
drug

>=65
years old

6 months

time in
database

# What is Needed in the CDM?

- **OMOP Vocabulary**
  to find the concepts

- **DRUG_EXPOSURE**
  to find individuals with exposure

- **OBSERVATION_PERIOD**
  to know people's time within the database

- **PERSON**
  to know year of birth

```
/*****************************************************************************
*       (Exercise 2) Warfarin New Users 65 or Older at Index
*****************************************************************************/

WITH CTE_DRUG_INDEX AS (
        SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
        FROM DRUG_EXPOSURE de
        WHERE de.DRUG_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
        )
        GROUP BY de.PERSON_ID
)
SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
        (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX,
        EXTRACT(YEAR FROM i.INDEX_DATE)-p.YEAR_OF_BIRTH AS AGE_AT_INDEX
FROM CTE_DRUG_INDEX i
        JOIN OBSERVATION_PERIOD op
                ON op.PERSON_ID = i.PERSON_ID
                AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
        JOIN PERSON p
                ON p.PERSON_ID = i.PERSON_ID
WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
AND EXTRACT(YEAR FROM i.INDEX_DATE)-p.YEAR_OF_BIRTH >= 65
ORDER BY i.PERSON_ID
```

```sql
/*************************************************************************
*        (Exercise 2) Warfarin New Users 65 or Older at Index
*************************************************************************/

WITH CTE_DRUG_INDEX AS (
        SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
        FROM DRUG_EXPOSURE de
        WHERE de.DRUG_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
        )
        GROUP BY de.PERSON_ID
)
SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
        (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX,
        EXTRACT(YEAR FROM i.INDEX_DATE)-p.YEAR_OF_BIRTH AS AGE_AT_INDEX
FROM CTE_DRUG_INDEX i
        JOIN OBSERVATION_PERIOD op
                ON op.PERSON_ID = i.PERSON_ID
                AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
        JOIN PERSON p
                ON p.PERSON_ID = i.PERSON_ID
WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
AND EXTRACT(YEAR FROM i.INDEX_DATE)-p.YEAR_OF_BIRTH >= 65
ORDER BY i.PERSON_ID
```

```
/************************************************************************
 *        (Exercise 2) Warfarin New Users 65 or Older at Index
 ************************************************************************/

WITH CTE_DRUG_INDEX AS (
        SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
        FROM DRUG_EXPOSURE de
        WHERE de.DRUG_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
        )
        GROUP BY de.PERSON_ID
)
SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
        (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX,
        EXTRACT(YEAR FROM i.INDEX_DATE)-p.YEAR_OF_BIRTH AS AGE_AT_INDEX
FROM CTE_DRUG_INDEX i
        JOIN OBSERVATION_PERIOD op
                ON op.PERSON_ID = i.PERSON_ID
                AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
        JOIN PERSON p
                ON p.PERSON_ID = i.PERSON_ID
WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
AND EXTRACT(YEAR FROM i.INDEX_DATE)-p.YEAR_OF_BIRTH >= 65
ORDER BY i.PERSON_ID
```

# How do I define new users of Warfarin with prior Atrial Fibrillation?

prior AFIB

index drug

time in database

6 months

# What is Needed in the CDM?

- **OMOP Vocabulary**
  to find the concepts

- **DRUG_EXPOSURE**
  to find individuals with exposure

- **OBSERVATION_PERIOD**
  to know people's time within the database

- **CONDITION_OCCURRENCE**
  to find presence of a disease

# Step 1: Start with the Ex 1 query

```sql
/******************************************************************************
*         (Exercise 3) Warfarin New Users With Prior AFIB
******************************************************************************/

WITH CTE_DRUG_INDEX AS (
        SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
        FROM DRUG_EXPOSURE de
        WHERE de.DRUG_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
        )
        GROUP BY de.PERSON_ID
),
CTE_DRUG_NEW_USERS AS (
        SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
                (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
        FROM CTE_DRUG_INDEX i
                JOIN OBSERVATION_PERIOD op
                        ON op.PERSON_ID = i.PERSON_ID
                        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
        WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
)
SELECT nu.*, MAX(nu.INDEX_DATE-co.CONDITION_START_DATE) AS DAYS_OF_CLOSEST_AFIB_PRIOR_TO_INDEX
FROM CTE_DRUG_NEW_USERS nu
        JOIN CONDITION_OCCURRENCE co
                ON co.PERSON_ID = nu.PERSON_ID
                AND co.CONDITION_START_DATE BETWEEN nu.OBSERVATION_PERIOD_START_DATE AND nu.OBSERVATION_PERIOD_END_DATE
WHERE co.CONDITION_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID =  313217 /*Atrial fibrillation*/
)
AND co.CONDITION_START_DATE < nu.INDEX_DATE
GROUP BY nu.PERSON_ID, nu.INDEX_DATE, nu.OBSERVATION_PERIOD_START_DATE, nu.OBSERVATION_PERIOD_END_DATE, nu.DAYS_BEFORE_INDEX
ORDER BY nu.PERSON_ID
```

# Step 2: Define Atrial Fibrillation

```
/*********************************************************************
*       (Exercise 3) Warfarin New Users With Prior AFIB
*********************************************************************/

WITH CTE_DRUG_INDEX AS (
        SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
        FROM DRUG_EXPOSURE de
        WHERE de.DRUG_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
        )
        GROUP BY de.PERSON_ID
),
CTE_DRUG_NEW_USERS AS (
        SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
                (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
        FROM CTE_DRUG_INDEX i
                JOIN OBSERVATION_PERIOD op
                        ON op.PERSON_ID = i.PERSON_ID
                        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
        WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
)
SELECT nu.*, MAX(nu.INDEX_DATE-co.CONDITION_START_DATE) AS DAYS_OF_CLOSEST_AFIB_PRIOR_TO_INDEX
FROM CTE_DRUG_NEW_USERS nu
        JOIN CONDITION_OCCURRENCE co
                ON co.PERSON_ID = nu.PERSON_ID
                AND co.CONDITION_START_DATE BETWEEN nu.OBSERVATION_PERIOD_START_DATE AND nu.OBSERVATION_PERIOD_END_DATE
WHERE co.CONDITION_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID =  313217 /*Atrial fibrillation*/
)
AND co.CONDITION_START_DATE < nu.INDEX_DATE
GROUP BY nu.PERSON_ID, nu.INDEX_DATE, nu.OBSERVATION_PERIOD_START_DATE, nu.OBSERVATION_PERIOD_END_DATE, nu.DAYS_BEFORE_INDEX
ORDER BY nu.PERSON_ID
```

# Step 3: Prior Atrial Fibrillation

```
/**********************************************************************
*         (Exercise 3) Warfarin New Users With Prior AFIB
**********************************************************************/

WITH CTE_DRUG_INDEX AS (
        SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
        FROM DRUG_EXPOSURE de
        WHERE de.DRUG_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
        )
        GROUP BY de.PERSON_ID
),
CTE_DRUG_NEW_USERS AS (
        SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
                (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
        FROM CTE_DRUG_INDEX i
                JOIN OBSERVATION_PERIOD op
                        ON op.PERSON_ID = i.PERSON_ID
                        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND
        WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
)
SELECT nu.*, MAX(nu.INDEX_DATE-co.CONDITION_START_DATE) AS DAYS_OF_CLOSEST_AFIB_PRIOR
FROM CTE_DRUG_NEW_USERS nu
        JOIN CONDITION_OCCURRENCE co
                ON co.PERSON_ID = nu.PERSON_ID
                AND co.CONDITION_START_DATE BETWEEN nu.OBSERVATION_PERIOD_START_DATE AND nu.OBSERVATION_PERIOD_END_DATE
WHERE co.CONDITION_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 313217 /*Atrial fibrillation*/
)
AND co.CONDITION_START_DATE < nu.INDEX_DATE
GROUP BY nu.PERSON_ID, nu.INDEX_DATE, nu.OBSERVATION_PERIOD_START_DATE, nu.OBSERVATION_PERIOD_END_DATE, nu.DAYS_BEFORE_INDEX
ORDER BY nu.PERSON_ID
```

Keeps condition within the same observable time, exclude if you want all time prior

# How do I define new users of Warfarin with prior Atrial Fibrillation?

prior AFIB

index drug

time in database

6 months

observation time

observation time

# New Users of Warfarin with prior Atrial Fibrillation

```sql
/********************************************************************************
*        (Exercise 3) Warfarin New Users With Prior AFIB
********************************************************************************/

WITH CTE_DRUG_INDEX AS (
        SELECT de.PERSON_ID, MIN(de.DRUG_EXPOSURE_START_DATE) AS INDEX_DATE
        FROM DRUG_EXPOSURE de
        WHERE de.DRUG_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID = 1310149 /*warfarin*/
        )
        GROUP BY de.PERSON_ID
),
CTE_DRUG_NEW_USERS AS (
        SELECT i.PERSON_ID, i.INDEX_DATE, op.OBSERVATION_PERIOD_START_DATE, op.OBSERVATION_PERIOD_END_DATE,
                (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) AS DAYS_BEFORE_INDEX
        FROM CTE_DRUG_INDEX i
                JOIN OBSERVATION_PERIOD op
                        ON op.PERSON_ID = i.PERSON_ID
                        AND i.INDEX_DATE BETWEEN op.OBSERVATION_PERIOD_START_DATE AND op.OBSERVATION_PERIOD_END_DATE
        WHERE (i.INDEX_DATE-op.OBSERVATION_PERIOD_START_DATE) >= 180
)
SELECT nu.*, MAX(nu.INDEX_DATE-co.CONDITION_START_DATE) AS DAYS_OF_CLOSEST_AFIB_PRIOR_TO_INDEX
FROM CTE_DRUG_NEW_USERS nu
        JOIN CONDITION_OCCURRENCE co
                ON co.PERSON_ID = nu.PERSON_ID
                AND co.CONDITION_START_DATE BETWEEN nu.OBSERVATION_PERIOD_START_DATE AND nu.OBSERVATION_PERIOD_END_DATE
WHERE co.CONDITION_CONCEPT_ID IN (
                SELECT DESCENDANT_CONCEPT_ID FROM CONCEPT_ANCESTOR WHERE ANCESTOR_CONCEPT_ID =  313217 /*Atrial fibrillation*/
)
AND co.CONDITION_START_DATE < nu.INDEX_DATE
GROUP BY nu.PERSON_ID, nu.INDEX_DATE, nu.OBSERVATION_PERIOD_START_DATE, nu.OBSERVATION_PERIOD_END_DATE, nu.DAYS_BEFORE_INDEX
ORDER BY nu.PERSON_ID
```

# Try on your own!

- Warfarin New Users 65 or Older at Index with Prior Atrial Fibrillation
   **8,207 individuals**


- Bonus:  Clipidogrel New Users 65 or Older at Index with Prior Atrial Fibrillation
   **3,148 individuals**

# Queries Can Be Automated

- Open up Google Chrome 

- Navigate to:
  http://localhost:8080/atlas/#/home

- In Atlas navigate to Cohorts 

- There should be a pre-existing cohort called "Warfarin New Users 65 or Older at Index with Prior Atrial Fibrillation."

# Queries Can Be Automated

# Break

## Please return in 15 minutes

# **Implementation**

Setup and Performing of an Extract Transform and Load process into the CDM

# Brief Review

- Foundational

  - OHDSI - Why and how

  - OMOP CDM - Standardizing structure & content

  - Real-world examples (SQL and Atlas)

# How do we create our own OMOP CDM instance?

**Extract**

**Transform**

**Load**

source$_1$

source$_2$

source$_3$

cdm

# ETL: Real world scenario

**Truven MarketScan Commercial Claims and Encounters (CCAE)**

**INPATIENT_SERVICES**

| enrolid | admdate | pdx | dx1 | dx2 | dx3 |
|---|---|---|---|---|---|
| 1570337021 | 5/31/2000 | 41071 | 41071 | 4241 | V5881 |

**Optum Extended SES**

**MEDICAL_CLAIMS**

| patid | fst_dt | diag1 | diag2 | diag3 | diag4 |
|---|---|---|---|---|---|
| 259000476532 | 5/30/2000 | 41071 | 27800 | 4019 | 2724 |

**Premier**

**PATICD_DIAG**

| pat_key | period | icd_code | icd_p |
|---|---|---|---|
| -17197140 | 1/1/2000 | 410.71 | P |
| -17197140 | 1/1/2000 | 414.01 | S |
| -17197140 | 1/1/2000 | 427.31 | S |
| -17197140 | 1/1/2000 | 496 | S |

**Japan Medical Data Center**

**DIAGNOSIS**

| member_id | admission_date | icd10_level4_code |
|---|---|---|
| M0041437 | 4/11/2013 | I214 |
| M0041437 | 4/11/2013 | A539 |
| M0041437 | 4/11/2013 | B182 |
| M0041437 | 4/11/2013 | E14- |

4 real observational databases, all containing an inpatient admission for a patient with a diagnosis of 'acute subendocardial infarction'
- Not a single table name the same…
- Not a single variable name the same….
- Different table structures (rows vs. columns)
- Different conventions (with and without decimal points)
- Different coding schemes (ICD9 vs. ICD10)

# What does it mean to ETL to OMOP CDM? Standardize **structure** and **content**

Truven MarketScan Commerical Claims and Encounters (CCAE)

**INPATIENT_SERVICES**

| enrolid | admdate | pdx | dx1 | dx2 | dx3 |
|---------|---------|-----|-----|-----|-----|
| 1570337021 | 5/31/2000 | 41071 | 41071 | 4241 | V5881 |

Structure optimized for large-scale analysis for clinical characterization, population-level estimation, and patient-level prediction

Truven MarketScan Commerical Claims and Encounters (CCAE)

**CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID |
|-----------|----------------------|------------------------|---------------------------|
| 157033702 | 5/31/2000 | 41071 | Inpatient claims - primary position |
| 157033702 | 5/31/2000 | 41071 | Inpatient claims - 1st position |
| 157033702 | 5/31/2000 | 4241 | Inpatient claims - 2nd position |
| 157033702 | 5/31/2000 | V5881 | Inpatient claims - 3rd position |

Content using international vocabulary standards that can be applied to any data source

Truven MarketScan Commerical Claims and Encounters (CCAE)

**CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID | CONDITION_SOURCE_CONCEPT_ID | CONDITION_CONCEPT_ID |
|-----------|----------------------|------------------------|---------------------------|-----------------------------|----------------------|
| 157033702 | 5/31/2000 | 41071 | Inpatient claims - primary position | 44825429 | 444406 |

# OMOP CDM = Standardized structure: same tables, same fields, same datatypes, same conventions across disparate sources

**Truven MarketScan Commerical Claims and Encounters (CCAE): INPATIENT_SERVICES**

| enrolid | admdate | pdx | dx1 | dx2 | dx3 |
|---|---|---|---|---|---|
| 1570337021 | 5/31/2000 | 41071 | 41071 | 4241 | V5881 |

**Optum Extended SES: MEDICAL_CLAIMS**

| patid | fst_dt | diag1 | diag2 | diag3 | diag4 |
|---|---|---|---|---|---|
| 259000476532 | 5/30/2000 | 41071 | 27800 | 4019 | 2724 |

**Premier: PATICD_DIAG**

| pat_key | period | icd_code | icd_pri_sec |
|---|---|---|---|
| -17197140 | 1/1/2000 | 410.71 | P |
| -17197140 | 1/1/2000 | 414.01 | S |
| -17197140 | 1/1/2000 | 427.31 | S |
| -17197140 | 1/1/2000 | 496 | S |

**JMDC: DIAGNOSIS**

| member_id | admission_date | icd10_level4_code |
|---|---|---|
| M0041437 | 4/11/2013 | I214 |
| M0041437 | 4/11/2013 | A539 |
| M0041437 | 4/11/2013 | B182 |
| M0041437 | 4/11/2013 | E14- |

**Truven CCAE: CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID |
|---|---|---|---|
| 157033702 | 5/31/2000 | 41071 | Inpatient claims - primary position |
| 157033702 | 5/31/2000 | 41071 | Inpatient claims - 1st position |
| 157033702 | 5/31/2000 | 4241 | Inpatient claims - 2nd position |
| 157033702 | 5/31/2000 | V5881 | Inpatient claims - 3rd position |

**Optum Extended SES: CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID |
|---|---|---|---|
| 259000474406532 | 5/30/2000 | 41071 | Inpatient claims - 1st position |
| 259000474406532 | 5/30/2000 | 27800 | Inpatient claims - 2nd position |
| 259000474406532 | 5/30/2000 | 4019 | Inpatient claims - 3rd position |
| 259000474406532 | 5/30/2000 | 2724 | Inpatient claims - 4th position |

**Premier : CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID |
|---|---|---|---|
| -171971409 | 1/1/2000 | 410.71 | Hospital record - primary |
| -171971409 | 1/1/2000 | 414.01 | Hospital record - secondary |
| -171971409 | 1/1/2000 | 427.31 | Hospital record - secondary |
| -171971409 | 1/1/2000 | 496 | Hospital record - secondary |

**JMDC : CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID |
|---|---|---|---|
| 4149337 | 4/11/2013 | I214 | Inpatient claims |
| 4149337 | 4/11/2013 | A539 | Inpatient claims |
| 4149337 | 4/11/2013 | B182 | Inpatient claims |
| 4149337 | 4/11/2013 | E14- | Inpatient claims |

- Consistent structure optimized for large-scale analysis
- Structure preserves all source content and provenance

# OMOP CDM = Standardized content: common vocabularies across disparate sources

**Truven CCAE: CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID | CONDITION_SOURCE_CONCEPT_ID | CONDITION_CONCEPT_ID |
|---|---|---|---|---|---|
| 157033702 | 5/31/2000 | 41071 | Inpatient claims - primary position | 44825429 | 444406 |

**Optum Extended SES: CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID | CONDITION_SOURCE_CONCEPT_ID | CONDITION_CONCEPT_ID |
|---|---|---|---|---|---|
| 259000474406532 | 5/30/2000 | 41071 | Inpatient claims - 1st position | 44825429 | 444406 |

**Premier : CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID | CONDITION_SOURCE_CONCEPT_ID | CONDITION_CONCEPT_ID |
|---|---|---|---|---|---|
| -171971409 | 1/1/2000 | 410.71 | Hospital record - primary | 44825429 | 444406 |

**JMDC : CONDITION_OCCURRENCE**

| PERSON_ID | CONDITION_START_DATE | CONDITION_SOURCE_VALUE | CONDITION_TYPE_CONCEPT_ID | CONDITION_SOURCE_CONCEPT_ID | CONDITION_CONCEPT_ID |
|---|---|---|---|---|---|
| 4149337 | 4/11/2013 | I214 | Inpatient claims | 45572081 | 444406 |

- Standardize source codes to be uniquely defined across all vocabularies
- No more worries about formatting or code overlap

- Standardize across vocabularies to a common referent standard (ICD9/10→SNOMED)
- Source codes mapped into each domain standard so that now you can talk across different languages

# ETL Process: Roles

- Members of the team

  –CDM expert

  –Local data expert

  –Data engineer

  –Person with medical knowledge

  –Business stakeholder

# ETL Process: Agile

# Example OHDSI ETL Process

## Analysis – Creation of ETL Specs/Stories

**Sprint 0**
- Location
- Care site
- Person
- Provider
- Condition
- Death
- Organization

**Sprint 1**
- Procedure Occurrence
- Observation
- Payer plan period
- Drug Cost
- Procedure Cost

**Sprint 2**
- Drug Exposure

**Sprint 3**
- Drug Era
- Condition Era
- Observation Period
- Visit Occurrence

**Sprint 4**
- Finalize ETL Specs

**For each table:**
- Backlog
- White Rabbit
- Vocabulary Mapping
- ETL specs

## Development – Implementation/Validation o

**Sprint 0**
- Initial Data On-boarding

**Sprint 1**
- Location
- Care site
- Person
- Provider
- Condition
- Death
- Organization

**Sprint 2**
- Procedure Occurrence
- Observation
- Payer plan period
- Drug Cost
- Procedure Cost

**Sprint 3**
- Drug Exposure

**Sprint 4**
- Drug Era
- Condition Era
- Observation Period
- Visit Occurrence

# OHDSI Resources for ETL

# Best Practices Documented

- [http://www.ohdsi.org/web/wiki/doku.php?id=documentation:etl_best_practices](http://www.ohdsi.org/web/wiki/doku.php?id=documentation:etl_best_practices)

# Getting WhiteRabbit

- [https://github.com/OHDSI/WhiteRabbit](https://github.com/OHDSI/WhiteRabbit)

- Click on "releases"



- "Latest Release" and download the WhiteRabbit zip file

# Getting WhiteRabbit

- Save the ZIP file somewhere and extract the files

- Double-click on the WhiteRabbit.jar to start the application.

# Working with WhiteRabbit

- Wiki: http://www.ohdsi.org/web/wiki/doku.php?id=documentation:software:whiterabbit

1. Specify the location of your data
   - Supports database connections as well as text files

2. Scanning your Database
   - Characterizes your data

# Specify the Location of Data

# Specify the Location of Data

# Scanning your Data

# Scanning your Data

# Scanning your Data

# Run the Scan Report on Your Data!

- Link on desktop


start_white_ra
bbit.sh

- Execute



**Execute File**                                                    − + ×

This text file 'start_white_rabbit.sh' seems to be an executable script. What do you want to do with it?

| ✓ Execute | Execute in Terminal | ⮝ Open | ✕ Cancel |

- WhiteRabbit appears

# Run the Scan Report on Your Data!

- Set the "Working Folder" to /home/ohdsi/whiterabbit/SynPUFSmall

- Press "Test connection"

- Move over to the "Scan" tab, and hit the "Add" button. Select the CSVs in the folder.

- Keep the default settings and press "Scan tables".

- Scan report is created in the folder you specified on the "Locations" tab as "ScanReport.xlsx".

# Reading the Scan

- Series of tabs in an XLSX file

  - **Overview Tab**
    provides the definition of each table analyzed, there will only be one tab of this type

  - **Table Tab(s)**
    a summary column for every column, there will be as many tabs as tables selected to analyze

# Overview Tab

- defines the tables you scanned

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | Table | Field | Type | Max lengt | N rows | N rows ch | Fraction empty | |
| 2 | Beneficiary_Summary.csv | DESYNPUF | varchar | 16 | -1 | 99999 | 0 | |
| 3 | Beneficiary_Summary.csv | BENE_BIRT | int | 8 | -1 | 99999 | 0 | |
| 4 | Beneficiary_Summary.csv | BENE_DEA | int | 8 | -1 | 99999 | 0.98449 | |
| 5 | Beneficiary_Summary.csv | BENE_SEX | int | 1 | -1 | 99999 | 0 | |
| 6 | Beneficiary_Summary.csv | BENE_RAC | int | 1 | -1 | 99999 | 0 | |
| 7 | Beneficiary_Summary.csv | BENE_ESR | varchar | 1 | -1 | 99999 | 0 | |
| 8 | Beneficiary_Summary.csv | SP_STATE | int | 2 | -1 | 99999 | 0 | |
| 9 | Beneficiary_Summary.csv | BENE_COL | int | 3 | -1 | 99999 | 0 | |
| 10 | Beneficiary_Summary.csv | BENE_HI_ | int | 2 | -1 | 99999 | 0 | |
| 11 | Beneficiary_Summary.csv | BENE_SMI | int | 2 | -1 | 99999 | 0 | |
| 12 | Beneficiary_Summary.csv | BENE_HM | int | 2 | -1 | 99999 | 0 | |
| 13 | Beneficiary_Summary.csv | PLAN_CVF | int | 2 | -1 | 99999 | 0 | |
| 14 | Beneficiary_Summary.csv | SP_ALZHD | int | 1 | -1 | 99999 | 0 | |
| 15 | Beneficiary_Summary.csv | SP_CHF | int | 1 | -1 | 99999 | 0 | |
| 16 | Beneficiary_Summary.csv | SP_CHRNI | int | 1 | -1 | 99999 | 0 | |
| 17 | Beneficiary_Summary.csv | SP_CNCR | int | 1 | -1 | 99999 | 0 | |
| 18 | Beneficiary_Summary.csv | SP_COPD | int | 1 | -1 | 99999 | 0 | |
| 19 | Beneficiary_Summary.csv | SP_DEPRE | int | 1 | -1 | 99999 | 0 | |
| 20 | Beneficiary_Summary.csv | SP_DIABE | int | 1 | -1 | 99999 | 0 | |
| 21 | Beneficiary_Summary.csv | SP_ISCHM | int | 1 | -1 | 99999 | 0 | |

# Table Tabs

- Definition of the Beneficiary_Summary.csv table and each record pertains to a synthetic medicare beneficiary

| # | Variable names | Labels |
|---|---|---|
| 1 | *DESYNPUF_ID* | DESYNPUF: Beneficiary Code |
| 2 | *BENE_BIRTH_DT* | DESYNPUF: Date of birth |
| 3 | *BENE_DEATH_DT* | DESYNPUF: Date of death |
| 4 | *BENE_SEX_IDENT_CD* | DESYNPUF: Sex |
| 5 | *BENE_RACE_CD* | DESYNPUF: Beneficiary Race Code |

**Variable Name: BENE_BIRTH_DT**

**Type:** Num

**Format:** YYYYMMDD

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | DESYNPUF_ID | Frequency | BENE_BIRTH_DT | Frequency | BENE_DEATH_DT | Frequency | BENE_SEX_IDENT_CD | Frequency |
| 2 | List truncated... | | 19421001 | 466 | | 98448 | 2 | 55211 |
| 3 | | | 19410401 | 434 | 20080901 | 151 | 1 | 44788 |
| 4 | | | 19390401 | 433 | 20080101 | 142 | | |
| 5 | | | 19431101 | 433 | 20081001 | 137 | | |
| 6 | | | 19400301 | 428 | 20080301 | 135 | | |
| 7 | | | 19410501 | 426 | 20081101 | 134 | | |
| 8 | | | 19390301 | 414 | 20080401 | 131 | | |
| 9 | | | 19400501 | 414 | 20080701 | 129 | | |
| 10 | | | 19410801 | 414 | 20080501 | 125 | | |
| 11 | | | 19391201 | 413 | 20080801 | 120 | | |
| 12 | | | 19431001 | 413 | 20081201 | 119 | | |
| 13 | | | 19411001 | 412 | 20080201 | 118 | | |
| | | | | 412 | 20080601 | 110 | | |
| | | | | 409 | | | | |

Beneficiary_Summary.csv

# Read the Scan Report

- Open up the scan report generated

- Go to the "Inpatient_Claims.csv" tab which represents claims processed from inpatient setting.

- What is the most common admitting diagnosis code?

- Hints:
  - ADMTNG_ICD9_DGNS_CD
  - ICD9 codes are in ###.## format
  - You can use ATLAS to look it up

# Rabbit In a Hat

- Already part of the WhiteRabbit download

- Uses the information from WhiteRabbit to help you produce documentation for the ETL process

- Helps you define the logic in a consistent way does not generate code to create ETL

# Getting Started

- Double-click on the RabbitInAHat.jar to start the application.

- File → Open Scan Report and navigate to the scan report that was just created.

# Process for Developing ETL

- Get the right people in the room

- Block off time necessary

- Map all the raw data tables to CDM tables

- Then go back through and provide detailed mapping information for each raw data table to CDM table

- Generate final ETL document

# Map Raw Tables to CDM Tables

# Map Raw Tables to CDM Tables

# Map Raw Tables
# to CDM Tables

# Map Raw Tables
# to CDM Tables

# Map Raw Tables
# to CDM Tables

# Map Raw Tables to CDM Tables

# Map Raw Tables to CDM Tables

Continue mapping raw tables to CDM tables until you feel confident you are bringing over as much raw data as necessary

# PERSON

- For today's example we'll start with the PERSON table

| Destination Field | Source Field | Logic | Comment |
|---|---|---|---|
| person_id | | | Autonumber |
| gender_concept_id | bene_sex_ident_cd | Source Value - Standard Concept Id<br>1 - 8507<br>2 - 8532<br><br>If gender is not 1 or 2, please discard person. | 1-Male<br>2-Female |
| year_of_birth | bene_birth_dt | Take first 4 digits (starting from left) | BENE_BIRTH_DT = YYYYMMDD |
| month_of_birth | bene_birth_dt | Take 5th and 6th digit starting from the left | BENE_BIRTH_DT = YYYYMMDD |
| day_of_birth | bene_birth_dt | Take last two digits starting from the left. | BENE_BIRTH_DT = YYYYMMDD |
| time_of_birth | | | N/A |
| race_concept_id | bene_race_cd | Source Value - Concept ID<br>1 - 8527<br>2 - 8516<br>3 - 0<br>5 - 0<br><br>Else set to 0. | 1-White<br>2-Black<br>3-Others<br>5-Hispanic |
| ethnicity_concept_id | bene_race_cd | Source Value - Concept ID<br>1 - 38003564<br>2 - 38003564<br>3 - 0<br>5 - 38003563<br><br>Else set to 0. | 1-White<br>2-Black<br>3-Others<br>5-Hispanic |
| location_id | sp_state_code<br>bene_county_cd | Use the BENE_COUNTY_CD and SP_STATE_CODE to lookup in the LOCATION table the LOCATION_ID. | |
| provider_id | | | N/A |
| care_site_id | | | N/A |
| person_source_value | desynpuf_id | | |
| gender_source_value | bene_sex_ident_cd | | |
| gender_source_concept_id | | | Set to 0. |
| race_source_value | bene_race_cd | | |
| race_source_concept_id | | | Set to 0. |
| ethnicity_source_value | bene_race_cd | | |
| ethnicity_source_concept_id | | | Set to 0. |

# DRUG_EXPOSURE

- Try drawing arrows from PRESCRIPTION_DRUG_EVENTS columns to DRUG_EXPOSURE columns



Focus on:

- – PERSON_ID
- – DRUG_EXPOSURE_START_DATE
- – QUANTITY
- – DAYS_SUPPLY
- – DRUG_SOURCE_VALUE

# DRUG_EXPOSURE

# DRUG_EXPOSURE

- Mapping source codes to standard terminology
  - Source to Source
  - Source to Standard

- Use standard query for both, just define filters needed

J9310 - *"Injection, rituximab, 100 mg"* → **Standard Query** → 46275076 - *"rituximab Injection"*

# Standard Query: Source to Standard

```sql
WITH CTE_VOCAB_MAP AS (
    SELECT c.concept_code AS SOURCE_CODE, c.concept_id AS SOURCE_CONCEPT_ID, c.concept_name AS
    SOURCE_CODE_DESCRIPTION, c.vocabulary_id AS SOURCE_VOCABULARY_ID, c.domain_id AS
    SOURCE_DOMAIN_ID, c.CONCEPT_CLASS_ID AS SOURCE_CONCEPT_CLASS_ID, c.VALID_START_DATE AS
    SOURCE_VALID_START_DATE, c.VALID_END_DATE AS SOURCE_VALID_END_DATE, c.INVALID_REASON AS
    SOURCE_INVALID_REASON,c1.concept_id AS TARGET_CONCEPT_ID, c1.concept_name AS
    TARGET_CONCEPT_NAME, c1.VOCABULARY_ID AS TARGET_VOCABUALRY_ID, c1.domain_id AS
    TARGET_DOMAIN_ID, c1.concept_class_id AS TARGET_CONCEPT_CLASS_ID, c1.INVALID_REASON AS
    TARGET_INVALID_REASON, c1.standard_concept AS TARGET_STANDARD_CONCEPT
    FROM CONCEPT C
        JOIN CONCEPT_RELATIONSHIP CR ON C.CONCEPT_ID = CR.CONCEPT_ID_1  AND CR.invalid_reason IS NULL
                AND cr.relationship_id = 'Maps to'
        JOIN CONCEPT C1 ON CR.CONCEPT_ID_2 = C1.CONCEPT_ID AND C1.INVALID_REASON IS NULL
    UNION
    SELECT source_code, SOURCE_CONCEPT_ID, SOURCE_CODE_DESCRIPTION, source_vocabulary_id, c1.domain_id
    AS SOURCE_DOMAIN_ID, c2.CONCEPT_CLASS_ID AS SOURCE_CONCEPT_CLASS_ID, c1.VALID_START_DATE AS
    SOURCE_VALID_START_DATE, c1.VALID_END_DATE AS SOURCE_VALID_END_DATE, stcm.INVALID_REASON AS
    SOURCE_INVALID_REASON,target_concept_id, c2.CONCEPT_NAME AS TARGET_CONCEPT_NAME,
    target_vocabulary_id, c2.domain_id AS TARGET_DOMAIN_ID, c2.concept_class_id AS TARGET_CONCEPT_CLASS_ID,
    c2.INVALID_REASON AS TARGET_INVALID_REASON, c2.standard_concept AS TARGET_STANDARD_CONCEPT
    FROM source_to_concept_map stcm
        LEFT OUTER JOIN CONCEPT c1 ON c1.concept_id = stcm.source_concept_id
        LEFT OUTER JOIN CONCEPT c2 ON c2.CONCEPT_ID = stcm.target_concept_id
    WHERE stcm.INVALID_REASON IS NULL
)
SELECT *
FROM CTE_VOCAB_MAP
/*EXAMPLE FILTERS*/
WHERE SOURCE_VOCABULARY_ID IN ('NDC')
AND TARGET_STANDARD_CONCEPT IN ('S')
```

# Standard Query: Source to Source

```sql
WITH CTE_VOCAB_MAP AS (
    SELECT c.concept_code AS SOURCE_CODE, c.concept_id AS SOURCE_CONCEPT_ID, c.CONCEPT_NAME AS
    SOURCE_CODE_DESCRIPTION, c.vocabulary_id AS SOURCE_VOCABULARY_ID, c.domain_id AS
    SOURCE_DOMAIN_ID, c.concept_class_id AS SOURCE_CONCEPT_CLASS_ID, c.VALID_START_DATE AS
    SOURCE_VALID_START_DATE, c.VALID_END_DATE AS SOURCE_VALID_END_DATE, c.invalid_reason AS
    SOURCE_INVALID_REASON, c.concept_ID as TARGET_CONCEPT_ID, c.concept_name AS
    TARGET_CONCEPT_NAME, c.vocabulary_id AS TARGET_VOCABULARY_ID, c.domain_id AS TARGET_DOMAIN_ID,
    c.concept_class_id AS TARGET_CONCEPT_CLASS_ID, c.INVALID_REASON AS
    TARGET_INVALID_REASON,c.STANDARD_CONCEPT AS TARGET_STANDARD_CONCEPT
    FROM CONCEPT c
    UNION
    SELECT source_code, SOURCE_CONCEPT_ID, SOURCE_CODE_DESCRIPTION, source_vocabulary_id, c1.domain_id
    AS SOURCE_DOMAIN_ID, c2.CONCEPT_CLASS_ID AS SOURCE_CONCEPT_CLASS_ID, c1.VALID_START_DATE AS
    SOURCE_VALID_START_DATE, c1.VALID_END_DATE AS SOURCE_VALID_END_DATE,stcm.INVALID_REASON AS
    SOURCE_INVALID_REASON,target_concept_id, c2.CONCEPT_NAME AS TARGET_CONCEPT_NAME,
    target_vocabulary_id, c2.domain_id AS TARGET_DOMAIN_ID, c2.concept_class_id AS TARGET_CONCEPT_CLASS_ID,
    c2.INVALID_REASON AS TARGET_INVALID_REASON, c2.standard_concept AS TARGET_STANDARD_CONCEPT
    FROM source_to_concept_map stcm
        LEFT OUTER JOIN CONCEPT c1 ON c1.concept_id = stcm.source_concept_id
        LEFT OUTER JOIN CONCEPT c2 ON c2.CONCEPT_ID = stcm.target_concept_id
    WHERE stcm.INVALID_REASON IS NULL
)
SELECT *
FROM CTE_VOCAB_MAP
/*EXAMPLE FILTERS*/
WHERE SOURCE_VOCABULARY_ID IN ('ICD9CM')
AND TARGET_VOCABULARY_ID IN ('ICD9CM')
```

# Example Filters: NDCs

- Source to Standard

```
WHERE SOURCE_VOCABULARY_ID IN ('NDC')
AND TARGET_STANDARD_CONCEPT IS NOT NULL
AND TARGET_INVALID_REASON IS NULL
AND DRUG_DATE BETWEEN SOURCE_VALID_START_DATE AND SOURCE_VALID_END_DATE
```

- Source to Source

```
WHERE SOURCE_VOCABULARY_ID IN ('NDC')
AND TARGET_VOCABULARY_ID IN ('NDC')
AND DRUG_DATE BETWEEN SOURCE_VALID_START_DATE AND SOURCE_VALID_END_DATE
```

| Some maps are date sensitive like NDC or DRGs | Review for incorrect mappings (e.g. source codes might map to multiple SOURCE_VOCAB_IDs) |
|---|---|

# Saving and Export to Document

- Save working document



- Export to document

# **Evaluation**

## Testing and Quality Assurance

# ACHILLES

- Interactive platform to visualize data in CDM
  - patient demographics
  - prevalence of conditions, drugs and procedures
  - distribution of values for clinical observations

- https://github.com/OHDSI/Achilles

# ETL Pitfalls

- Privacy Issues
  - Removal of ICD9/10 codes that are considered privacy issues, such as death or sexual abuse
  - Using "fake" date in Death table to indicate a death

- Patient Cleansing
  - Test patients

- Differing Business Rules
  - Institutions decide not to follow vocabulary classifications

# Conclusion

# Join the journey

Interested in OHDSI?

Join the Journey:

http://www.ohdsi.org/join-the-journey/

Questions:

http://forums.ohdsi.org/

APPENDIX



Lasers

# USAGI

- Tool to help in mapping codes from a source system into the standard terminologies stored in OMOP Vocabulary

http://www.ohdsi.org/web/wiki/doku.php?id=documentation:software:usagi

# USAGI Exercise

Sample File

# USAGI Exercise

# USAGI Exercise

# USAGI Exercise

# USAGI Exercise

# USAGI Exercise

**source_to_concept_map**

| | |
|---|---|
| **source_code** | varchar(50) |
| **source_concept_id** | int |
| **source_vocabulary_id** | varchar(20) |
| **source_code_description** | varchar(255) |
| **target_concept_id** | int |
| **target_vocabulary_id** | varchar(20) |
| **valid_start_date** | date |
| **valid_end_date** | date |
| **invalid_reason** | varchar(1) |

| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| | source_code | source_concept_id | source_vocabulary_id | source_code_description | target_concept_id | target_vocabulary_id | valid_start_date | valid_end_date | invalid_reason |
| | Negative | 0 | categorical_lab_map | Negative | 45878583 | LOINC | 1/1/70 | 12/31/99 | |
| | Colorless | 0 | categorical_lab_map | Colorless | 45880448 | LOINC | 1/1/70 | 12/31/99 | |
| | + | 0 | categorical_lab_map | Positive | 45884084 | LOINC | 1/1/70 | 12/31/99 | |
| | Non-React | 0 | categorical_lab_map | Non-React | 4305306 | SNOMED | 1/1/70 | 12/31/99 | |
| | Normal | 0 | categorical_lab_map | Normal | 45884153 | LOINC | 1/1/70 | 12/31/99 | |
| | Not Detected | 0 | categorical_lab_map | Not Detected | 45880296 | LOINC | 1/1/70 | 12/31/99 | |
| | Positive | 0 | categorical_lab_map | Positive | 45884084 | LOINC | 1/1/70 | 12/31/99 | |