# CDM Workshop, Part 1

Common Data Model Working Group

2022-03-08
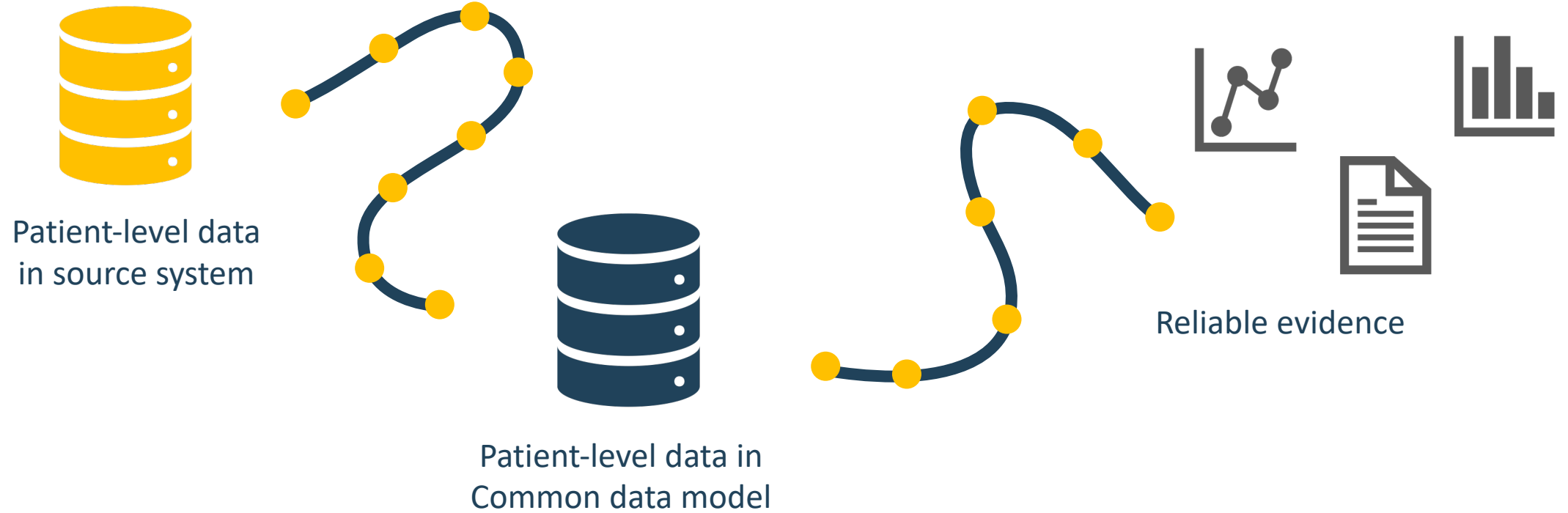
# The Journey from Data to ...

# The Journey from Data to …
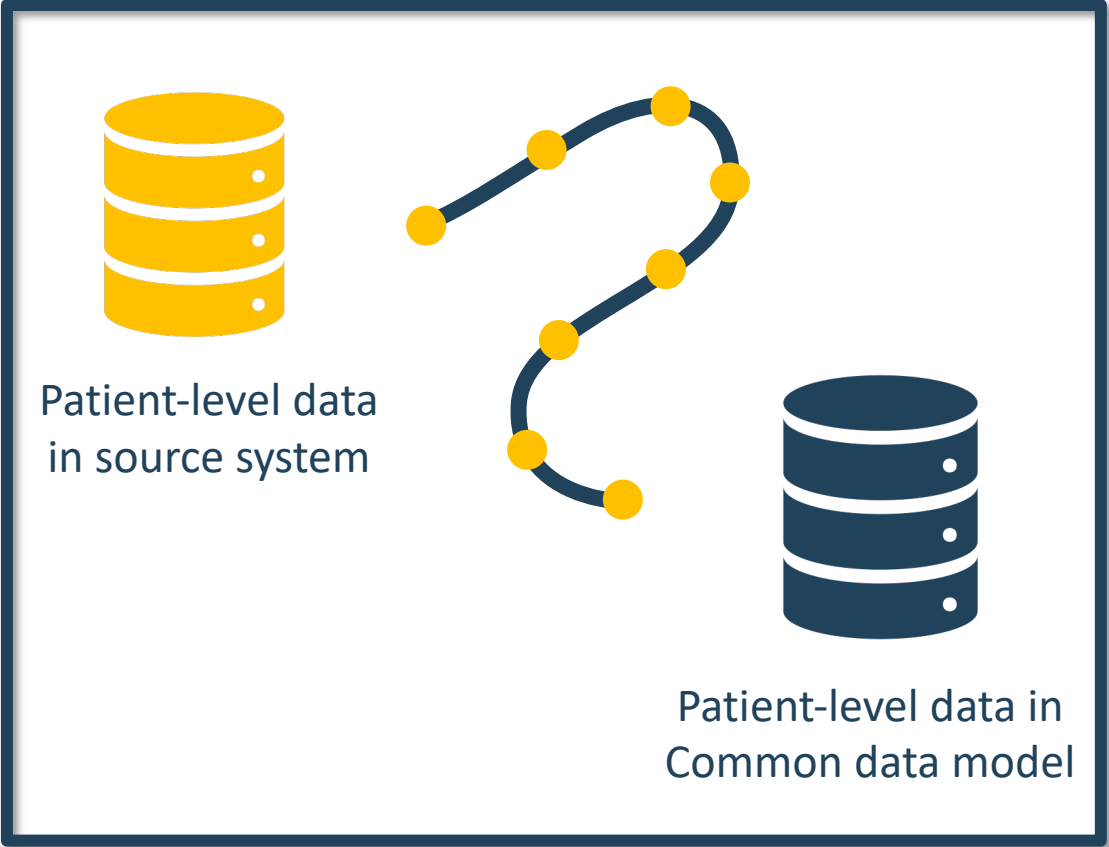
# Data

# The Journey from Data to Reliable Evidence

Patient-level data
in source system

Patient-level data in
Common data model

Reliable evidence

# The Journey from Source to Standardized Data



Patient-level data in source system

Patient-level data in Common data model

Reliable evidence

# The Journey from Source to Standardized Data



Technical Considerations for Setup

Data Governance

White Rabbit

Rabbit In A Hat

Usagi

Data Quality Dashboard

ACHILLES

Technical Considerations for Studies

Patient-level data in source system

Patient-level data in Common data model

# The Journey from Source to Standardized Data

Technical Considerations for Setup

Data Governance

White Rabbit

Rabbit In A Hat

Usagi

Data Quality Dashboard

ACHILLES

Technical Considerations for Studies

Patient-level data
in source system

Patient-level data in
Common data model

# The Journey from Source to Standardized Data

**Technical Considerations for Setup**

**Data Governance**

**White Rabbit**

**Rabbit In A Hat**

Frank DeFalco

Kristin Kostka

Maxim Moinat

# Technical Considerations for Setup

Frank DeFalco

CDM Workshop 2022

2022-03-08

# Relational Database Requirements

- Start with a compatible relational database platform
  - Microsoft SQL Server
  - Oracle
  - PostgreSQL
  - Amazon RedShift
  - Impala
  - IBM Netezza
  - Google BigQuery
  - Microsoft PDW
  - Apache Spark
  - SQLite

# Setting up your CDM

- Choose a currently supported CDM Version (5.3, **5.4**)

  - CommonDataModel
    - R Package to generate a CDM on your database platform
    - https://github.com/OHDSI/CommonDataModel

    ```
    CommonDataModel::executeDdl(
        connectionDetails = cd,
        cdmVersion = "5.4",
        cdmDatabaseSchema = "ohdsi_journey"
    )
    ```

  - A vocabulary schema is required to exist and contain a vocabulary content distribution along with any instance of the CDM
  - Athena
    - https://athena.ohdsi.org/
    - Web site where you can download vocabulary content

# CDM v5.4 Schema



Before you start loading your data, consider more than the technical...

# What is data governance?

- A process to formally outline how organizational data will be managed and controlled

- Ensures that data is consistent and trustworthy and is not misused

# Why does it matter in the OMOP CDM?

- OMOP CDM is a person-centric model

- The model can retain attributes that may be considered personal identified information (PII) or protected health information (PHI)

- There are many ways a site may treat their OMOP CDM to uphold their governance protocols

# The OHDSI data holder's responsibility

- In OHDSI research, it is the responsibility of each data holder to:
  - **know,**
  - **understand,**
  - **and follow**
- local data governance processes related to use of your OMOP CDM instance.

# Sites have different rules about patient data

- ***In Europe:***
- General Data Protection Regulation (GDPR) is a regulation in EU law on data protection and privacy in the European Union (EU) and the European Economic Area (EEA)
  - *Implemented 25 May 2018*

- ***In the US:***
- The Health Insurance Portability and Accountability Act (HIPAA) of 1996 establishes national standards to protect individuals' medical records and other personal health information





**The HIPAA Privacy Rule**

The HIPAA Privacy Rule establishes national standards to protect individuals' medical records and other personal health information and applies to health plans, health care clearinghouses, and those health care providers that conduct certain health care transactions electronically. The Rule requires appropriate safeguards to protect the privacy of personal health information, and sets limits and conditions on the uses and disclosures that may be made of such information without patient authorization. The Rule also gives patients rights over their health information, including rights to examine and obtain a copy of their health records, and to request corrections.

The Privacy Rule is located at 45 CFR Part 160 and Subparts A and E of Part 164.

Click here to view the combined regulation text of all HIPAA Administrative Simplification Regulations found at 45 CFR 160, 162, and 164.

# Where do privacy concerns arise?

- Privacy issues usually happen around:
- Date of Birth
- Patient location
- Precise clinical event date
- Psychological or mentally related clinical conditions
- Death date
- Death cause (usually non-medical related death, murder, car accident etc.)

# Where should I be checking my OMOP CDM?

- Date Fields Across Domains
- *_source_value
- Any string fields

- Extra sensitive tables: PERSON, LOCATION, PROVIDER, OBSERVATION, NOTE and NOTE_NLP

# More detail on the OMOP CDM Wiki

# Alphabet soup: governance edition

Institutional Review Board (IRB)

Data use agreement (DUA)

Data transfer agreements (DTA)

Data use request (DUR)

CITI training
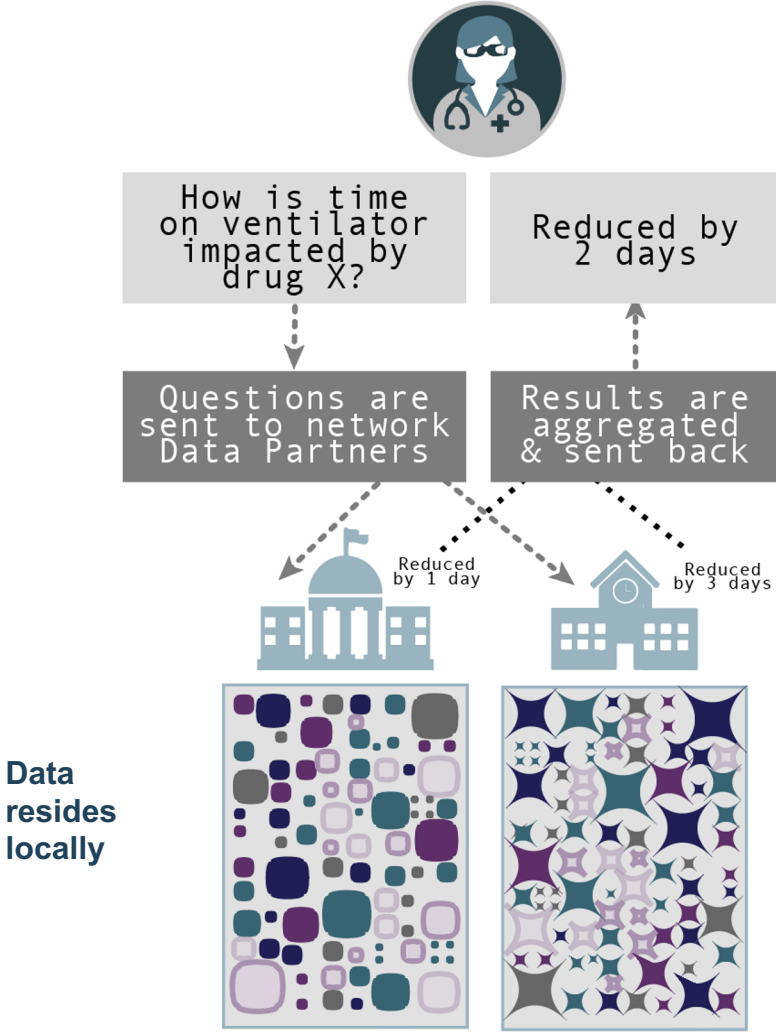
Data stewards
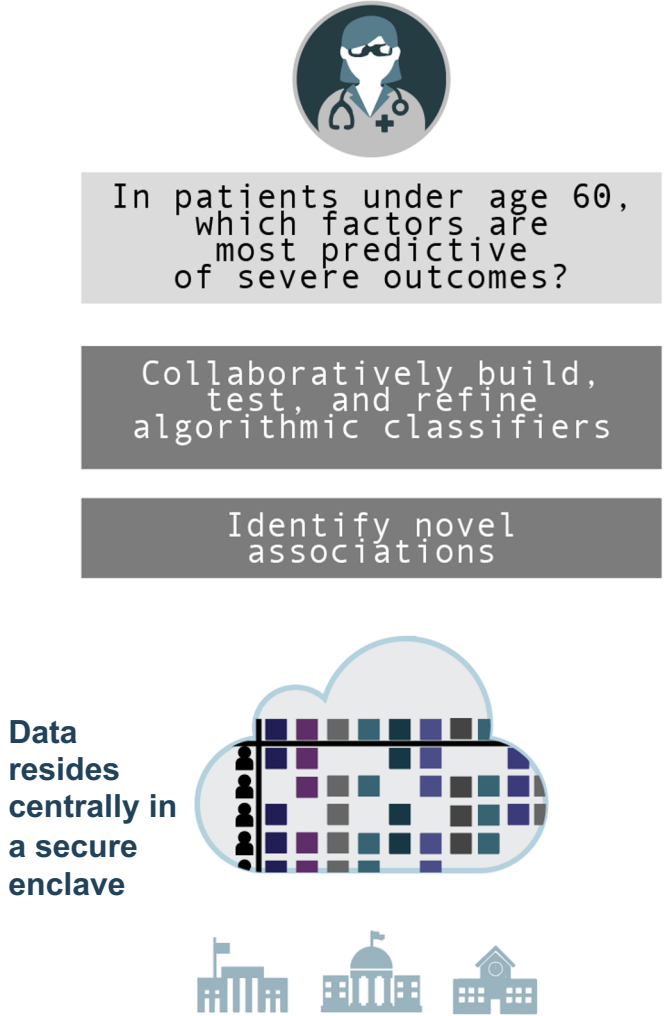
Limited data set (LDS)

Synthetic data set

# Types of data querying

## Federated querying

How is time on ventilator impacted by drug X?

Reduced by 2 days

Questions are sent to network Data Partners

Results are aggregated & sent back

Reduced by 1 day

Reduced by 3 days

**Data resides locally**

## Centralized analytics

In patients under age 60, which factors are most predictive of severe outcomes?

Collaboratively build, test, and refine algorithmic classifiers

Identify novel associations

**Data resides centrally in a secure enclave**

Credit: N3C

# Closing thoughts

- Governance is not one-size-fits-all

- Consult your local IRB and privacy office for guidance

- When in doubt, ask for help before you break any rules!

Google's secret cache of medical data includes names and full details of millions - whistleblo...

HELLO
my name is
That Guy

**TECH**

**Congre...ional Democrats demand details on Google's use of patient data by Dec. 6**

# ETL Process



Data experts and CDM experts together design the ETL

People with medical knowledge create the code mappings

ETL Documentation

ETL

All are involved in quality control

A technical person implements the ETL

OHDSI Tools

White Rabbit

Rabbit In a Hat

Usagi

White Rabbit

ACHILLES

DQD

Rabbit In a Hat

# WR - Overview

- **Create Scan Report**
  – Overview of tables, fields and values
  – Record counts and statistics

- **Support for:**
  – Delimited text (csv)
  – Sas7bdat files
  – Eight RDBMS

- **Fake data generator**
  – Based on Scan Report
  – Random values per column

# Why create a scan report?

- Helps understanding the structure of your source data
- Input for Rabbit in a Hat

**Table/Field Overview**

| Table | Field | Type | Max length | N rows | N rows ch | Fraction emp |
|-------|-------|------|-----------|--------|-----------|--------------|
| allergies | start | date | 10 | 619 | 619 | 0 |
| allergies | stop | date | 10 | 619 | 619 | 0.904685 |
| allergies | patient | character | 36 | 619 | 619 | 0 |
| allergies | encounter | character | 36 | 619 | 619 | 0 |
| allergies | code | character | 9 | 619 | 619 | 0 |
| allergies | description | character | 24 | 619 | 619 | 0 |
| | | | | | | |
| careplans | id | character | 36 | 2939 | 2939 | 0 |
| careplans | start | date | 10 | 2939 | 2939 | 0 |
| careplans | stop | date | 10 | 2939 | 2939 | 0.380061 |
| careplans | patient | character | 36 | 2939 | 2939 | 0 |
| careplans | encounter | character | 36 | 2939 | 2939 | 0 |
| careplans | code | character | 15 | 2939 | 2939 | 0 |
| careplans | description | character | 62 | 2939 | 2939 | 0 |
| careplans | reason_co | character | 14 | 2939 | 2939 | 0.090507 |
| careplans | reason_de | character | 69 | 2939 | 2939 | 0.090507 |
| | | | | | | |
| conditions | start | date | 10 | 7898 | 7898 | 0 |
| conditions | stop | date | 10 | 7898 | 7898 | 0.458091 |
| conditions | patient | character | 36 | 7898 | 7898 | 0 |
| conditions | encounter | character | 36 | 7898 | 7898 | 0 |
| conditions | code | character | 7 | 7898 | 7898 | 0.545455 |
| conditions | description | character | 80 | 7898 | 7898 | 0 |
| | | | | | | |
| encounter | id | character | 36 | 34275 | 34275 | 0 |
| encounter | start | date | 10 | 34275 | 34275 | 0 |
| encounter | stop | date | 10 | 34275 | 34275 | 0 |

Tabs: Overview | allergies | careplans | conditions | encounte

**Value counts**

| marital | Frequency | race | Frequency | ethnicity | Frequency | gender | Frequency | birthplace | Frequency |
|---------|-----------|------|-----------|-----------|-----------|--------|-----------|-----------|-----------|
| M | 622 | white | 846 | irish | 235 | M | 572 | Boston | 142 |
| | 344 | hispanic | 112 | italian | 145 | F | 558 | Springfiel | 30 |
| S | 166 | black | 82 | english | 102 | | 2 | Worceste | 28 |
| | | asian | 70 | puerto_ric | 72 | | | Lowell | 22 |
| | | native | 20 | french | 72 | | | Brockton | 21 |
| | | other | 1 | german | 64 | | | Cambridg | 18 |
| | | Unknown | 1 | chinese | 51 | | | Methuen | 18 |
| | | | | polish | 49 | | | Newton | 17 |
| | | | | american | 39 | | | Quincy | 16 |
| | | | | portugues | 37 | | | Framingha | 16 |
| | | | | french_ca | 35 | | | Lynn | 12 |
| | | | | african | 33 | | | Arlington | 12 |
| | | | | west_indi | 28 | | | Weymout | 12 |
| | | | | dominica | 21 | | | New Bedf | 12 |
| | | | | american_ | 20 | | | Lawrence | 11 |
| | | | | russian | 20 | | | Haverhill | 11 |
| | | | | scottish | 19 | | | Fitchburg | 11 |
| | | | | asian_indi | 19 | | | Marshfiel | 10 |
| | | | | mexican | 18 | | | Somervill | 10 |
| | | | | swedish | 17 | | | Barnstable | 10 |
| | | | | central_ar | 13 | | | Fall River | 9 |
| | | | | greek | 12 | | | Waltham | 9 |

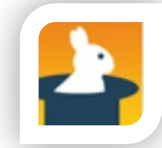Tabs: immunizations | medications | observations | organizations | **patients** | proce ...
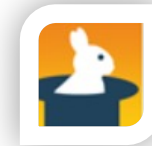
# RiaH - Overview

- Facilitates mapping workshop
- Source and Target tables
  - Source tables read from WR Scan Report
  - Target tables from an OMOP CDM definition (or custom)
- Interactively document transformation rules by linking:
  - Table to Table
  - Field to Field

# Rabbit in a Hat makes it easy to explore your source data structure!

# Create the mapping specification

## Table to Table

Tables

| Source | CDMV5.4 |
|--------|---------|
| patients.csv | person |

## Field to Field

| Source | CDMV5.4 |
|--------|---------|
| id | *person_id |
| gender | person_source_value |
| birthdate | *gender_concept_id |
| | gender_source_value |
| | *year_of_birth |
| | month_of_birth |
| | day_of_birth |
| | birth_datetime |

## Document mapping decisions and logic

Details

General information

Source: patients.csv.birthdate
Target: person.year_of_birth

Logic

Take Year from birthdate
Drop rows with year < 1900

# RiaH – Export mapping definition

Word document

Markdown documents

Html

Condition_occurrence
Drug_exposure
Measurement
Observation
Observation_period
Person
Procedure_occurrence
Visit Occurrence

# Person

**Contents**

Person

## Person

### Reading from Synthea table patients.csv

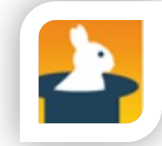| Destination Field | Source field | Logic | Comment field |
|---|---|---|---|
| person_id | | | |
| gender_concept_id | gender | When gender = 'M' then set gender_concept_id to 8507, when gender = 'F' then set to 8532 | Drop any rows with missing/unknown gender. |
| year_of_birth | birthdate | Take year from birthdate | |
| month_of_birth | birthdate | Take month from birthdate | |
| day_of_birth | birthdate | Take day from birthdate | |
| birth_datetime | birthdate | With midnight as time 00:00:00 | |
| death_datetime | deathdate | With midnight as time 00:00:00 | |
| race_concept_id | race | When race = 'WHITE' then set as 8527, when race = 'BLACK' then set as 8516, when race = 'ASIAN' then set as 8515, otherwise set as 0 | |
| ethnicity_concept_id | race ethnicity | When race = 'HISPANIC', or when ethnicity in ('CENTRAL_AMERICAN', 'DOMINICAN', 'MEXICAN', 'PUERTO_RICAN', 'SOUTH_AMERICAN' ) then set as 38003563, otherwise set as 0 When race = 'HISPANIC', or when ethnicity in ('CENTRAL_AMERICAN', 'DOMINICAN', 'MEXICAN', 'PUERTO_RICAN', 'SOUTH_AMERICAN' ) then set as 38003563, otherwise set as 0 | |

Publish your ETL document!
Example: https://ohdsi.github.io/ETL-Synthea/articles/person.html

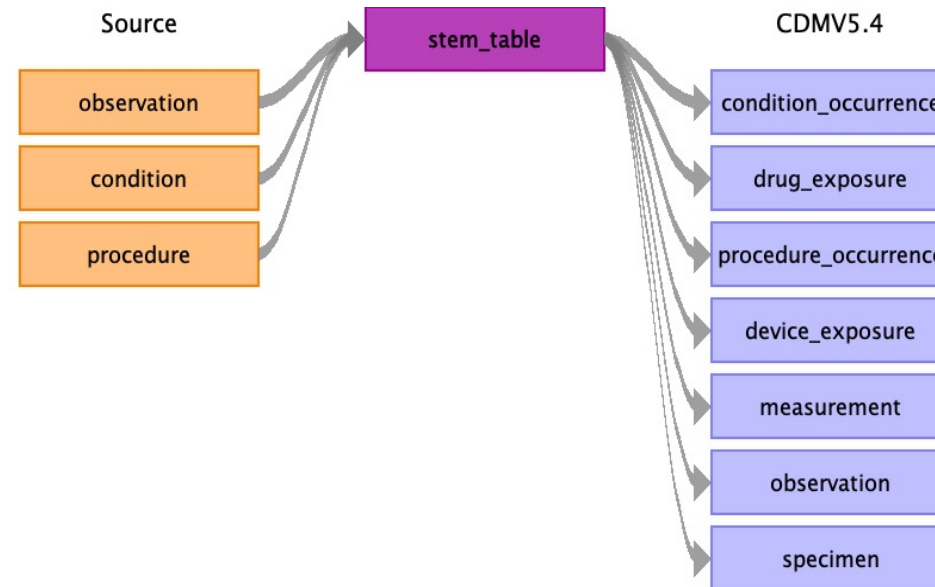# Vocabulary drives data movement:
# the RiaH Stem table

e.g. not all codes from condition table will be actual condition occurrences.

- SNOMED:161833006 "Abnormal weight gain " is an Observation

The RiaH STEM Table is a domain-agnostic intermediate.

Relations to OMOP domains are provided.

# Mapping Workshop Checklist

✓ General database information

✓ Technical specifications

✓ Data Provenance

✓ Use Cases

✓ Goals of the transformation

✓ Refreshes (data & ETL)

✓ Data restrictions (ethical)

✓ Contact for access (data manager)

💡 Get to know the source data!

# Final remark

Rabbit-in-a-Hat is only for creating the **structural mapping** specification and documenting mapping decisions.

"Which source fields are the input for which OMOP CDM fields."

⚠️ The mapping of values to standard concepts, semantic mapping, is a separate mapping process.
More on semantic mapping with Usagi next week!

# ETL Process