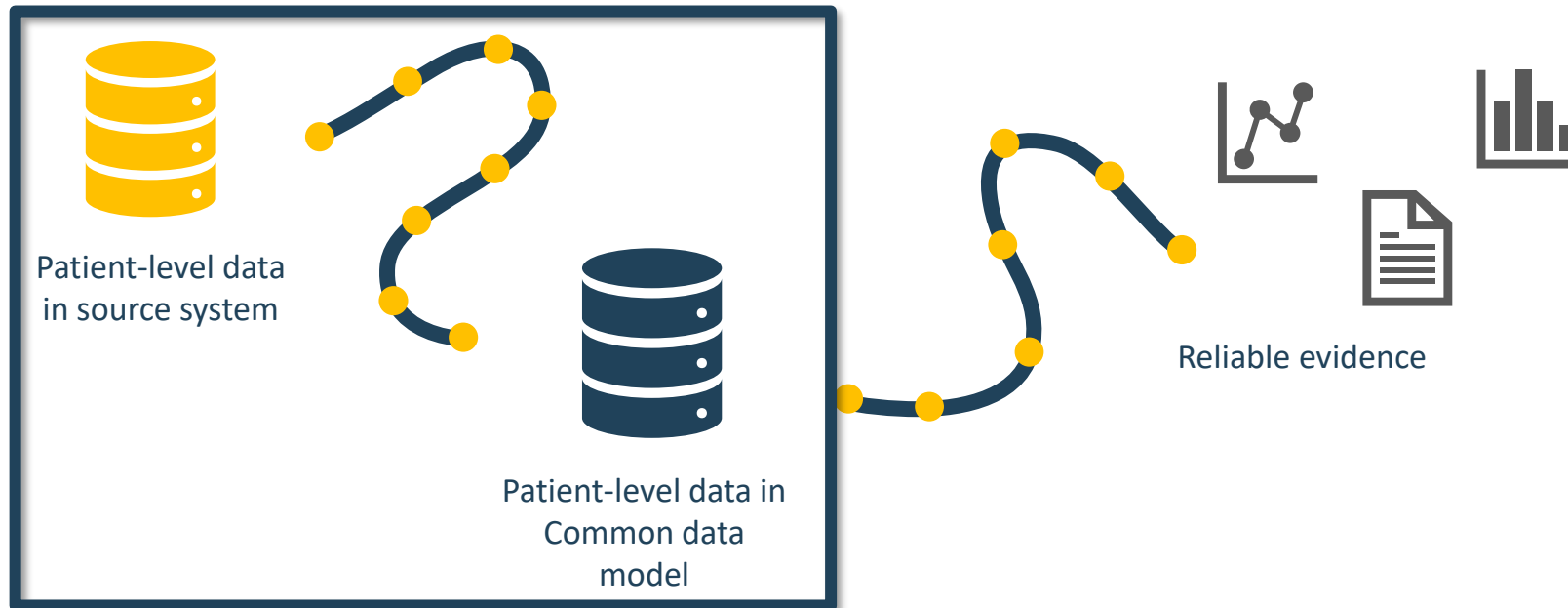# OMOP Conversion Process

# ETL
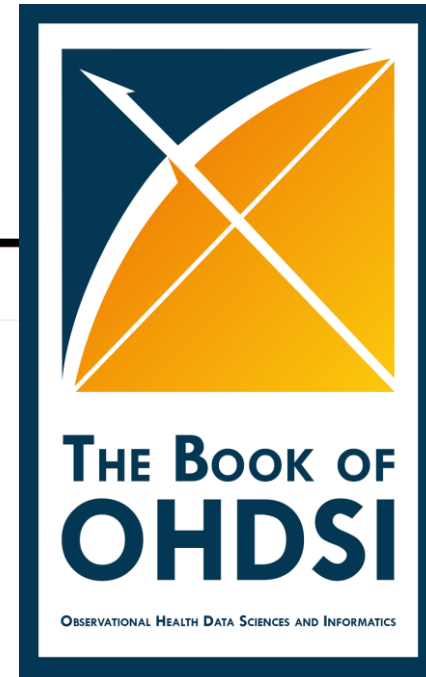
- Extract Transform Load

- In order to get from our native/raw data into the OMOP CDM we need to design and develop and ETL process



Patient-level data in source system

Patient-level data in Common data model

Reliable evidence

- Goal in ETLing is to standardize the format and terminology

# ETL Process

# ETL Process



Data experts and CDM experts together design the ETL

People with medical knowledge create the code mappings

ETL Documentation

ETL

A technical person implements the ETL

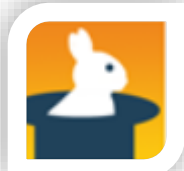All are involved in quality control
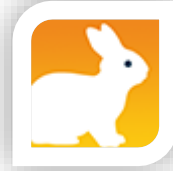
OHDSI Tools

White Rabbit | Rabbit In a Hat | Usagi | White Rabbit | ACHILLES | DQD | Rabbit In a Hat

# Designing the ETL



Data experts and CDM experts together design the ETL

People with medical knowledge create the code mappings

ETL Documentation

ETL

A technical person implements the ETL

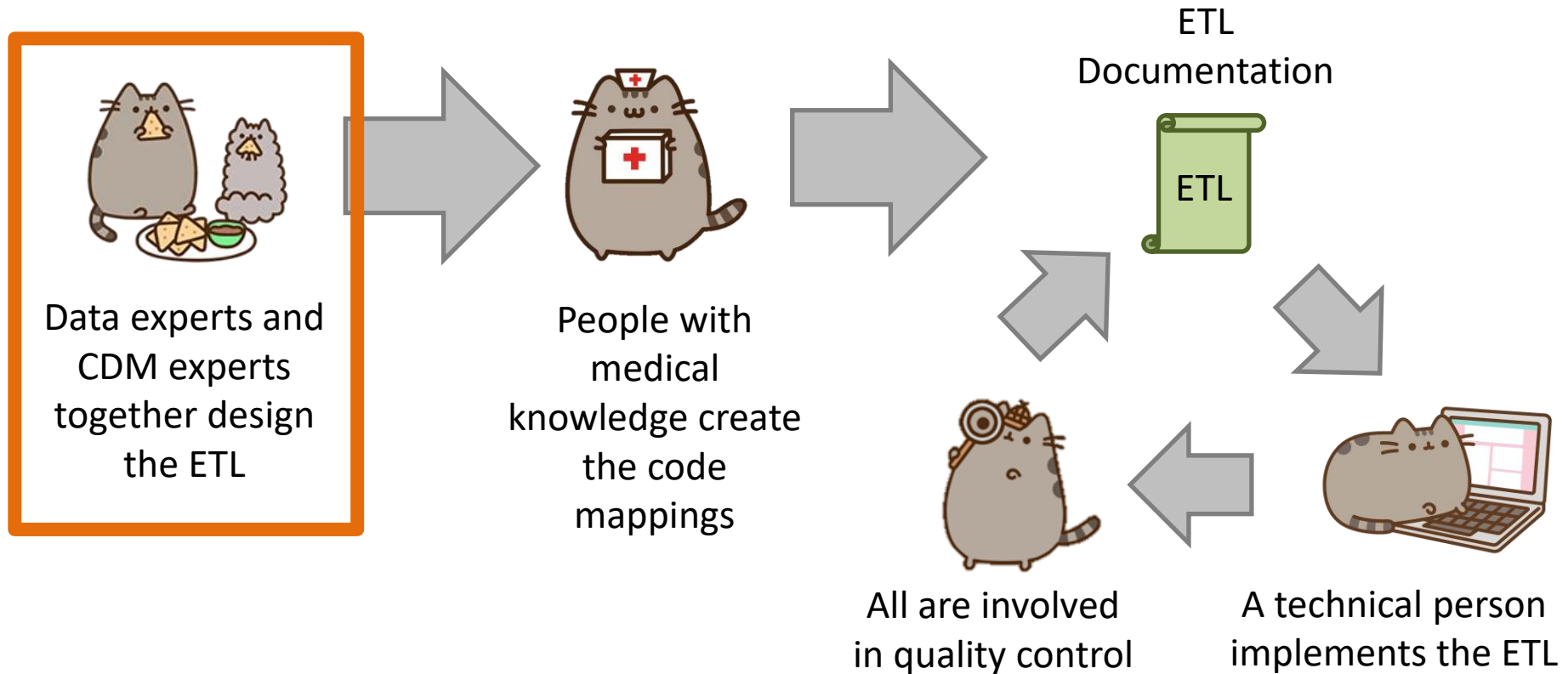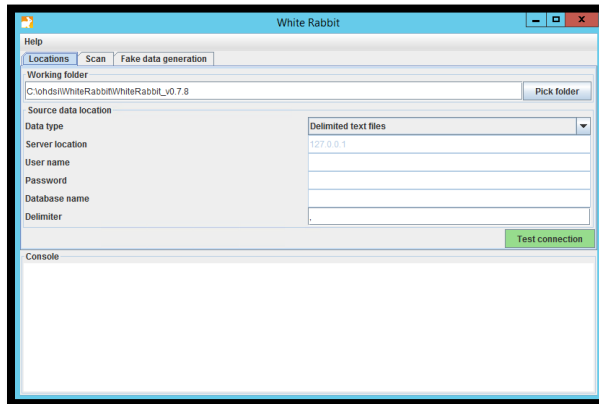All are involved in quality control

# White Rabbit

- White Rabbit scans source data & creates a csv report on the source data

- The scan can be used to:
  - Learn about your source data
  - Help design the ETL
  - Used by Rabbit In a Hat

# WR Output – ScanReport.xlsx

## Table/Field Overview

| Table | Field | Description | Type | Max length | N rows |
|---|---|---|---|---|---|
| pop | der_sex | | character | 1 | 16374539 |
| pop | der_yob | | double pre | 6 | 16374539 |
| pop | pat_id | | character | 64 | 16374539 |
| pop | pat_hash_id | | character | 16 | 16374539 |
| pop | pmtx_flag | | numeric | 1 | 16374539 |
| pop | anon_ims_pat_id | | character | 11 | 16374539 |
| pop | pat_region | | character | 2 | 16374539 |
| pop | pat_state | | character | 2 | 16374539 |
| pop | pat_zip3 | | character | 3 | 16374539 |
| pop | grp_indv_cd | | character | 1 | 16374539 |
| pop | mh_cd | | character | 1 | 16374539 |
| pop | enr_rel | | character | 2 | 16374539 |
| pop | temp_col1 | | character | 0 | 16374539 |
| pop | temp_col2 | | character | 0 | 16374539 |
| pop | load_row_id | | bigint | 9 | 16374539 |
| | | | | | |
| claims_diag_lk | person_source_valu | | character | 64 | 2992046684 |
| claims_diag_lk | event_start_date | | date | 10 | 2992046684 |
| claims_diag_lk | event_end_date | | date | 10 | 2992046684 |

## Value counts

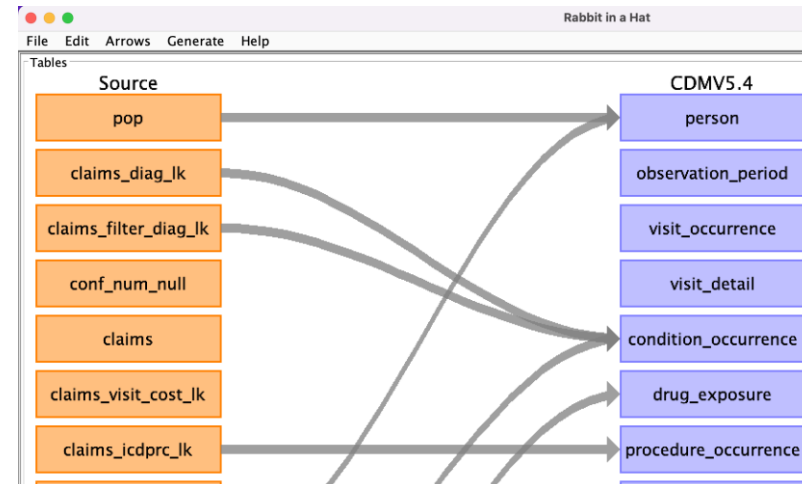| | A | B | C | D | |
|---|---|---|---|---|---|
| 1 | der_sex | Frequency | der_yob | Frequency | pa |
| 2 | F | 50479 | 1991.0 | 2030 | Li |
| 3 | M | 49514 | 1992.0 | 1970 | |
| 4 | U | 7 | 1990.0 | 1947 | |
| 5 | | | 1989.0 | 1908 | |
| 6 | | | 1988.0 | 1873 | |
| 7 | | | 1994.0 | 1872 | |
| 8 | | | 1995.0 | 1806 | |
| 9 | | | 1993.0 | 1805 | |
| 10 | | | 1996.0 | 1716 | |
| 11 | | | 1986.0 | 1676 | |
| 12 | | | 1987.0 | 1643 | |
| 13 | | | 1985.0 | 1633 | |
| 14 | | | 1983.0 | 1588 | |
| 15 | | | 1981.0 | 1581 | |
| 16 | | | 1984.0 | 1576 | |
| 17 | | | 1970.0 | 1555 | |
| 18 | | | 1980.0 | 1553 | |

pop | claims_diag_lk | claims

# Rabbit in a Hat



- Read and display a White Rabbit scan document

- Provides a graphical interface to allow a user to connect source data to CDM tables

# RiaH - Output

Word document

Markdown documents

Html

# Vocabulary Mapping



Data experts and CDM experts together design the ETL

People with medical knowledge create the code mappings

ETL Documentation

ETL

A technical person implements the ETL

All are involved in quality control

# Usagi

- When the Vocabulary does not contain your source terms you will need to create a map to OMOP Vocabulary Concepts

- Usagi helps you to:
  - Find best matches, automatically and/or manually
  - Automatic matching based on text similarities (itf/df)
  - Create 'source to concept map'

# Overview - Steps

1. Get a copy of the Vocabulary from ATHENA
2. Download Usagi
3. **Have Usagi build an index on the Vocabulary**

One-time setup

4. Load your source codes and let Usagi process them
5. Review and update suggested mappings with someone who has medical knowledge
6. Export codes into the SOURCE_TO_CONCEPT_MAP

# Implementing the ETL

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

# ETL Implementation

There are multiple tools available to implement your ETL



Your choice will largely depend on the size and complexity of the ETL design. And the tools available to you.

# ETL Implementation

## General Flow of Implementation

person

observation_period

visit_occurrence

A good rule of thumb is to always create the PERSON table first

The VISIT_OCCURRENCE table must be created before the standardized clinical data tables as they all refer to the VISIT_OCCURRENCE_ID

condition_occurrence

observation

drug_exposure

procedure_occurrence

measurement

Additional clinical data tables...

# Quality Control



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

# Quality

What tools are available to check that the CDM logic was implemented correctly?

Rabbit-in-a-Hat Test Case Framework

Achilles

DataQualityDashboard (DQD)

# Unit Test Cases

- Testing your CDM builder is important:
  - ETL is often complex, increasing the danger of making mistakes that go unnoticed

  - CDM can update

  - Source data structure/contents can change over time

- Rabbit-In-a-Hat can construct unit tests, or small pieces of code that can automatically check single aspects of the ETL design

# Achilles

Achilles is a data characterization and quality tool available for download here:

https://github.com/OHDSI/Achilles

For an example of how it was run for some sample data, that R script is located here:

https://github.com/OHDSI/Tutorial-ETL/blob/master/materials/Achilles/achillesRun.R

# DataQualityDashboard (DQD)

- Runs a prespecified set of data quality checks and thresholds on the CDM

**DATA QUALITY ASSESSMENT**

**SYNTHEA SYNTHETIC HEALTH DATABASE**

Results generated at 2019-08-22 14:15:06 in 29 mins

OVERVIEW
METADATA
RESULTS
ABOUT

| | Verification | | | | Validation | | | | Total | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Pass | Fail | Total | % Pass | Pass | Fail | Total | % Pass | Pass | Fail | Total | % Pass |
| Plausibility | 159 | 21 | 180 | 88% | 283 | 0 | 283 | 100% | 442 | 21 | 463 | 95% |
| Conformance | 637 | 34 | 671 | 95% | 104 | 0 | 104 | 100% | 741 | 34 | 775 | 96% |
| Completeness | 369 | 17 | 386 | 96% | 5 | 10 | 15 | 33% | 374 | 27 | 401 | 93% |
| Total | 1165 | 72 | 1237 | 94% | 392 | 10 | 402 | 98% | 1557 | 82 | 1639 | **95%** |

# Common ETL Issues

**Non-standard Vocabulary**
Codes mapped to OMOP vocabulary aren't mapped to a 'Standard'

**Multiple Input on Records**
Some records will contain multiple coding systems and text. A hierarchy must be selected to avoid duplicate records

**Non-Clinical Events**
Due to text options in EHR Data, many options are not clinical events (e.g. 'Tuesday' or 'XXYZ'). These records will be scrubbed to ensure quality of data converted to OMOP.

**Multiple records for one concept mapping**
Picking one of the multiple standard vocabulary mapping to create the OMOP CDM record instead of one record per mapping

**Abnormal values**
Unconventional values in data asset (i.e. Negative or 0 as value_as_number)

**Incorrect logic - Observation_Period**
Observation_Period table populated incorrectly. Observation period does not cover the entire period of time where events are recorded for a person

**Wrong type_concept_id**
Use of the wrong type_concept_id or misunderstanding the definition of this field

**Missing CDM tables**
OMOP CDM tables missing due to misunderstanding on how to populate the table.
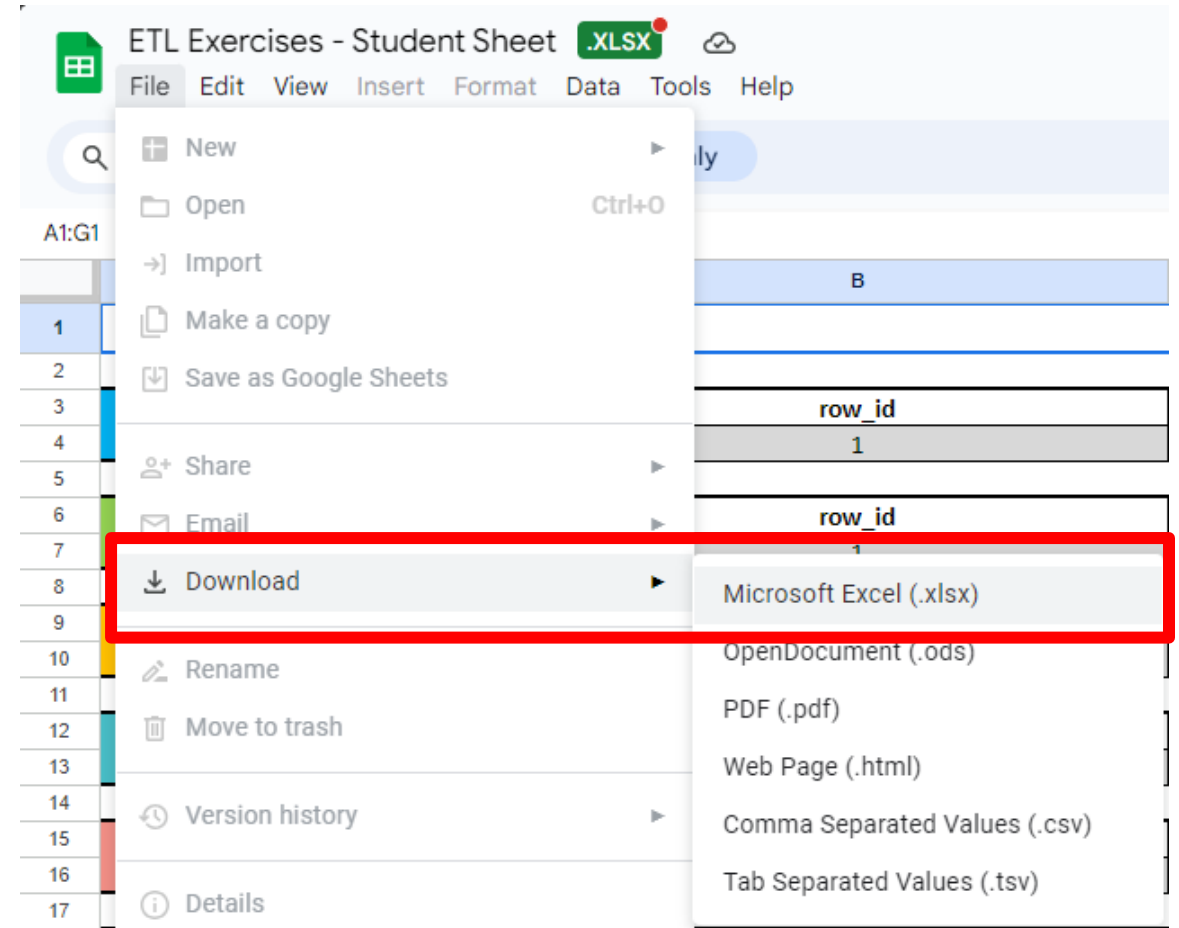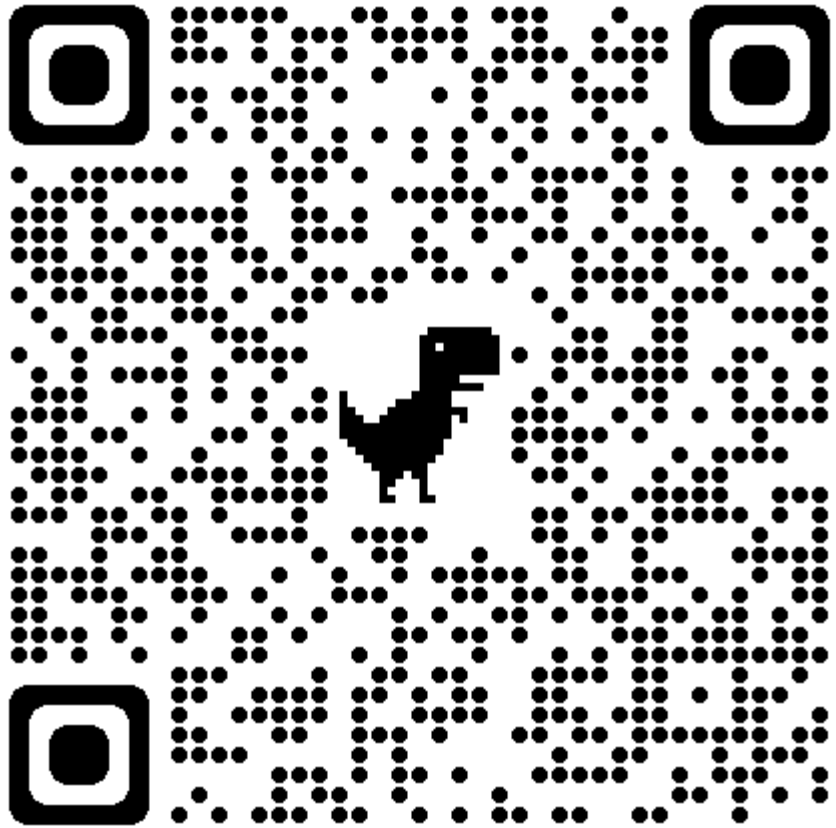
**Incorrect logic - Visit_Occurrence**
Visit_Occurrence table populated incorrectly

# Exercise Instructions

- Download a copy of the exercises at:

# Exercise Instructions

- Using the native data provided, map it to the OMOP CDM using the template provided in the *ETL Development_1000* sheet

- If you have spare time, do the same for the *ETL Development_1005* and *ETL Development_1010* sheets

# Thank you!

Mui Van Zandt     mui.vanzandt@iqvia.com

Gyeol Song     gyeol.song@iqvia.com