

CDM Workshop, Part 2

Common Data Model Working Group 2022-03-15



The Journey from Data to ...



The Journey from Data to ...

Data



The Journey from Data to Reliable Evidence



Patient-level data in Common data model









in source system

Patient-level data in Common data model





7



Vocabulary Mapping & Usagi



Melanie Philofsky

Data Quality



ACHILLES



Putting it all Together



Anthony Molinaro

Frank DeFalco

Clair Blacketer



Vocabulary Mapping and USAGI

Melanie Philofsky CDM Workshop 2022 2022-03-15



OMOP Vocabulary

- Who, what, when, why & where?
- How does it all work?
- Custom semantic mapping
 - Usagi
 - Source to Concept Map table
 - 2 Billionaires
- Maintenance



OMOP Vocabulary tables

- Where they came from
- Who does all this work
- What they are
 - Quick review of source code to concept_id mapping
- How do I get these
- When & why



Where and who

- Public domain
- Proprietary sources
- Homegrown OMOP concepts

- OHDSI Vocabulary team
- Vocabulary maintenance
 - Source release cycles
 - Mapping to standards
- OHDSI Vocabulary GitHub page



When & why do you need the OMOP Vocabularies

- Now's a good time to download them
- Can't properly populate a CDM without them
- Can't use OHDSI tools without them



Athena.ohdsi.org

NΔ	SEARCH	DOWNLOAD	LOGIN	0			
Search							
aspirin			Sea	arch			
 Usage of quotation marks forces an exact-match search In case of a typo, or if there is a similar spelling of the word, the most similar result will be presented 							
Explore domains							
Drugs 5,250,974	Conditions	Conditions Procedures ^{698,822} [®] 1733,500					
Devices	Observations	Measure	ements				



Vocabulary Version

ATHENA - OHDSI VOCABULARIES REPOSITORY



@ 2015-2022, Odysseus Data Services, Inc. All rights reserved

Version 1.12.2.8.210316.0857

OMOP Vocabulary version: v5.0 04-FEB-22

Report application issue

Report vocabulary content issues



OMOP Standardized Vocabulary Tables

- Concept
- Concept Relationship
- Concept Ancestor
- Source to Concept Map

- Synonym
- Relationship
- Concept Class
- Domain
- Drug Strength
- Vocabulary



Structure of OMOP Vocabulary



All content: concepts in concept

Direct relationships between concepts in concept_relationship Multi-step hierarchical relationships pre-processed into concept_ancestor





The Source for Source Codes

May come from international terminology or code system

SNOMED, ISBT



May come from a country specific terminology or code system

Read, BDPM, ICD10CN, CVX



May be free text strings

• Centimeter, Intravenous, Cigarette Smoker



May come from a source specific code system

• EHRs, CRFs, Registries, etc.



Mapping Source Codes to Concept_IDs

Concept_ID	 Map to standard or non-standard concept_id Do not map to classification concept_id*
Source Code Vocabulary	 Vocabulary is known: use to retrieve concept_id Vocabulary is unknown, but domain is known: use domain to retrieve concept_id
Custom/Local Source Code	 Custom or local source code: custom map to standard concept_id

*Classification concept_ids are used outside the clinical event tables for research



Source Code Mapping – Scenario 1

Scenario

 Source code comes from an OMOP supported Vocabulary

Solution

 Use using the following condition to perform the mapping: Where <source code> = CONCEPT.concept_code and <source vocabulary> = CONCEPT.vocabulary_id

Source code	Source vocabulary	Code description	CONCEPT.concept_id
61462000	SNOMED	Malaria	438067
A663D00	Read	Zika Fever	45489770
A92.3	ICD10CN	West Nile Virus Infection	1404276



Source Codes Mapping – Scenario 2

Scenario

• Source code is a text string

Solution

 Use using the following condition to perform the mapping: Where <source string> = CONCEPT.concept_name and <source domain> = CONCEPT.domain_id

Source string	Source domain	Source table/field	CONCEPT.concept_id
Centimeter	Unit	Vital Signs/ unit for height measurement	8582
Intravenous	Route	Drug/ route for drug administration	4171047
Male	Gender	Demographics	8507



Source Codes Mapping – Scenario 3

Scenario

 Source data does not map to an OHDSI supported vocabulary

Solution

- Ask OHDSI
- Create custom mapping using one of the following two methods:
 - Source to Concept Map
 - 2 Billionaires



Custom Semantic Mapping

- Usagi
- Source to Concept Map table
- Concept & Concept Relationship tables
- When do you NOT create a custom concept_id?





🕤 Usagi - tes	ticpcCodesMap	ping.csv										_				x
<u>File Edit V</u>	liew <u>H</u> elp															
Status	Source code	e Source term	Frequency	CodeText	Match score	Concept ID	Concept nam	ne Domain	Concept cl	ass Vocabulary	Concept code S	tandard con	Parents	Children	Commer	t
Approved	K87.00	Hypertension.	694195	Hypertensie	0.81	316866	Hypertensive.	Condition	Clinical Fin	di SNOMED	38341003 S		1	27		-
Approved	L99.00	Other diseas.	680422	Andere ziekte.	0.47	0	Unmapped	Oradition	Oliniaal Fin		400044000		0	0	Too generic	
Approved	D01.00	Abdominal p.,	675917	Gegeneralis	0.01	197988	Generalized	Condition	Clinical Fin	di SNOMED	102014000 S		1	102		
Inchecked	T86.00	Hypothyroidi	667283	HypothyreoÄ	1.00	4113642	Hypothyroidi	Condition	Clinical Fin	di SNOMED	286910004 S		c 1	0		-
Source code								••••••								00000
	Sou	urce code			So	urce term				Frequency			(CodeText		
S99.00				Skin disease	other			675817				Andere ziekte	(n) huid/subc	utis		
Target conc	epts															
Con	cept ID	Concept n	ame	Domain		Concept class	Vo	ocabulary	Con	cept code	Standard conce	ept	Parents		Children	
4317258	l	Disorder of skin	Con	dition	Clinica	I Finding	SNOMED		95320005	5	5	2		193		
															Remove conce	pt
Search														Lin		
Query									Filters							
														O dia secti		
Use sou	rce term as qu	erv							E Fill	er by user select	led concepts	Filter by co	oncept class:	Z-dig nonbli	ii code	
0.0000									Filt	er standard con	cepts	Filter by vo	cabulary:	ABMS	-	
U Query:									✓ Inc	lude source term	ns	Filter by do	main:	Condition	-	
Results														-		
Sco	re	Term	Concept	ID	Concept n	ame	Domai	n Co	ncept class	Vocabulary	Concept co	de Standar	d concept	Parents	Children	Π
0.75	Skin	disease	4317258	Disorde	rofskin		Condition	Clinica	al Finding	SNOMED	95320005	S	2		193	-
0.65	Skin	Disease, Funga	137213	Dermal	mycosis		Condition	Clinica	al Finding	SNOMED	14560005	S	3		12	
0.57	AIDS	with skin dise	4224566	Skin dis	order associat	ed with AIDS	Condition	Clinica	al Finding	SNOMED	421394009	S	2		2	=
0.56	Chro	nic skin disease	9 4134132	Chronic	disease of ski	n	Condition	Clinica	al Finding	SNOMED	128236002	S	2		26	
0.55	Dise	ase, Otologic	378161	Disorde	r of ear		Condition	Clinica	al Finding	SNOMED	25906001	S	4		43	
0.55	Dise	ase, Hers	4163346	Glycoge	n storage dise	ase, type VI	Condition	Clinica	al Finding	SNOMED	29291001	S	2		0	
0.55	Othe	r peripheral va	321052	Periphe	ral vascular dis	ease	Condition	Clinica	al Finding	SNOMED	400047006	S	1		44	
0.55	Othe	r peripheral va	4119612	Lower li	mb ischemia		Condition	Clinica	al Finding	SNOMED	233961000	S	2		3	
0.55	Dise	ase, Ormond	4176725	Retrope	ritoneal fibrosi	s	Condition	Clinica	al Finding	SNOMED	49120005	S	1		3	
0.54	Path	ological fractur	73571	Patholo	gical fracture		Condition	Clinica	al Finding	SNOMED	268029009	S	1		21	
0.52	Dise	ase, Tooth	4122115	Tooth di	sorder		Condition	Clinica	al Finding	SNOMED	234947003	S	3		58	
0.52	Dise	ase, Lip	135858	Disorde	roflip		Condition	Clinica	al Finding	SNOMED	90678009	S	3		35	-
0.51	Diea	aca Olliar	4113600	Multiple	concenital exc	etneie	Condition	Clinics	Einding	SNOMED	254044004	9	6		0	1
													Repla	ce concept	Add concep	it.
-																
omment:															Appr	ove
naroyed / tot	al: 5/1037 9	8% of total free	Mency													



Source_to_Concept_Map

Example source code = Pediatric interventional cardiologist

Field	Source_cod e	Source_c oncept_i d	Source_vo cabulary_i d	Source_code _description	Target_c oncept_ id	Target_v ocabular y_id	Valid_start_ date	Valid_end_ date	Invalid _reaso n
Required Field?	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	No
Value	Pediatric interventio nal cardiologist	Unique identifier >2 billion	prov_speci alty	Pediatric intervention al cardiologist	903276	Medicare Specialty	01/01/1970	12/31/2099	





- Create a Concept
- Create the Concept Relationship





Example source code = Pediatric interventional cardiologist

Concept_id	Concept_ name	Domain _id	Vocabula ry_id	Concept _class_id	Standard _concept	Concept_ code	Valid_sta rt_date	Valid_en d_date	Invalid_r eason
Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	No
Unique id > 2 billion 21000000	Pediatric interventi onal cardiologi st	Provider	prov_spe cialty	Physician Specialty		Pediatric interventi onal cardiologi st	01/01/19 70	12/31/20 99	



Concept Relationship

Example source code = Pediatric interventional cardiologist

Field	Concept_id_1	Concept_id_2	Relationship_id	Valid_start_da te	Valid_end_dat e	Invalid_reason
Required Field?	Yes	Yes	Yes	Yes	Yes	No
Record #1 value	2100000000	903276	Maps to	01/01/1970	12/31/2099	
Record #2 value	903276	2100000000	Mapped from	01/01/1970	12/31/2099	



Comparison of both methods

STCM

Concept & Concept Relationship

- Designed for insert statements
- Need to create an additional SQL statement to do a lookup on the STCM during the ETL

- Need to modify OHDSI provided tables
 - Concept, Concept Relationship, Vocabulary tables
- Enriches your Vocabulary
- Enables the display of source concepts in Atlas
- Great for internal use cases



When do you NOT create a custom concept?

- New source code, already supported vocabulary
- New vocabulary, strong use case



Maintenance

- Every time you update your OHDSI vocabularies, re-run your ETL
- Complete an analysis of your CDM
 - Review top unmapped values
 - Review for duplicates
 - Are source codes now represented by an OHDSI supported concept_id?
 - Update STCM or Concept & Concept Relationship tables



Resources

- Athena This is a web browser for the most up to date vocabularies
 - <u>http://Athena.ohdsi.org</u>
- USAGI: Download the program, request enhancements, raise issues
 - <u>https://github.com/OHDSI/Usagi</u>
- OHDSI/OMOP Vocabulary GitHub: Log issues/defects & request enhancements
 - <u>https://github.com/OHDSI/Vocabulary-v5.0</u>
- Forums: Ask questions, open discussions, raise ideas
 - https://forums.ohdsi.org
- CDM GitHub page: Log issues/defects & request enhancements
 - <u>https://github.com/OHDSI/CommonDataModel</u>
- CDM wiki page: All conventions,
 - <u>https://ohdsi.github.io/CommonDataModel/faq.html</u>
- Book of OHDSI: Central repository for OHDSI knowledge
 - <u>https://ohdsi.github.io/TheBookOfOhdsi/</u>



Clair Blacketer CDM Workshop 2022 2022-03-15



ETL Process





Data experts and CDM experts together design the ETL People with medical knowledge create the code mappings



ETL

Documentation



All are involved in quality control

A technical person implements the ETL





Data Quality – Thoughts From the FDA

Semicerrenss 2000, it emier 2000)

E. Quality Assurance (QA) and Quality Control (QC)

Asse: Use i Investigators should fully understand the *quality assurance (QA)* and *quality control (QC)* procedures used by the data holders and how these procedures could have an effect on the integrity of the data and the overall validity of the study. FDA recommends that investigators address the following topics:

Assesting and Relevance

The strength of RWE submitted in support of a regu depends on the clinical study methodology and the r accrual and data quality control (data assurance)) ar underlying data. In general, FDA does not endorse o

- The frequency and type of any data error corrections or changes in data adjudication policies implemented by the data holders during the relevant period of data collection;
- A description of any peer-reviewed publications examining data quality and/or validity, including the relationships of the investigators with the data source(s);
- Any updates and changes in coding practices (e.g., ICD codes) across the study period that are relevant to the outcomes of interest;
- Any changes in key data elements during the study time frame and their potential effect on the study; and
- A report on the extent of missing data over time (i.e., the percentage of data not available for a particular variable of interest) and a discussion on the procedures (e.g., exclusion, imputation) employed to handle this issue. Investigators should also address the implications of the extent of missing data on study findings and the missing data methods used.

Framework for FDA's Real-World Evidence Program



Data Quality – Thoughts From the EMA

Establish a certification process for data sources '

'Data quality is not a static construct and is context, disease and question dependent and dependent on the healthcare system. Assessments need to be constant and documented every time the data is refreshed'

Public consultation comments

In order to include novel data sources as evidence sources for regulatory decision-making, it is critical to understand how much the regulators can rely on the data. Thus, a capability to characterise the quality of data is a strategic objective for regulators. While pre-defining quality is challenging as need is often driven by the question, it is possible to define some generalised elements for which quality could be defined.

What this means for stakeholders:

A data quality framework will support the trust of patients and healthcare professionals in the decisions reached by regulators when Big Data underpins those decisions. It will aid the choice of data source selected for a study (including those by industry) and it will inform the assessment of the study results and the benefit-risk dossier by regulators.



What is Data Quality?





Data Quality Check Types

Check Type	Check Description
Person Completeness	The number and percent of persons in a database that do not have a least one record in the <i>CDM table</i> .
Is Required	The number and percent of records with a NULL value in a <i>CDM field</i> of a <i>CDM table</i> that is considered not nullable.
Is Primary Key	The number and percent of records that have a duplicate value in the <i>CDM field</i> of the <i>CDM table</i> .
Is Foreign Key	The number and percent of records that have a value in a <i>lookup field</i> of a <i>CDM table</i> that does not exist in the <i>lookup table</i> .
Concept Domain	The number and percent of records that have a concept_id that does not conform to the <i>domain</i> .
Is Standard Valid Concept	The number and percent of records that do not have a standard, valid concept in the <i>CDM field</i> of a <i>CDM table</i> .



Data Quality Check Types

Check Type	Check Description
Standard Concept Completeness	The number and percent of records with a value of 0 in the standard concept field <i>CDM field</i> in the <i>CDM table</i> .
Plausible Temporal After	The number and percent of records with a value in a <i>CDM field</i> of a <i>CDM table</i> that occurs prior to a <i>plausible date</i> .
Plausible Value Low	For a given <i>concept_id</i> and <i>unit_concept_id</i> pair, the number and percent of records with a value lower than the <i>plausible low value</i> .
Plausible Gender	For a given <i>concept_id</i> , the number and percent of records associated with persons with an <i>implausible gender</i> .



Data Quality Check Types

	Verification	Validation
Plausibility	6	1
Conformance	7	1
Completeness	4	1

20 Check *Types*



Data Quality Check Totals

	Verification	Validation
Plausibility	1855	287
Conformance	563	80
Completeness	327	12

Total 3,124 Checks





RESULTS Show IBM MARKETSCAN COMMERCIAL < **CLAIMS AND ENCOUNTERS** DATABASE ____ Ŧ **OVERVIEW** Ŧ METADATA Ŧ ABOUT + Ŧ _____

IBM MARKETSCAN COMMERCIAL CLAIMS AND ENCOUNTERS DATABASE

Results generated at 2019-09-06 22:20:12 in 7 hours

						Column	visibility CSV
Show	TATUS	CONTEXT	CATEGORY	Search:	% RECORDS		
Ð	PASS	Verification	Completeness	None	FIELD	The number and percent of records with a NULL value in the range_high of the MEASUREMENT. (Threshold=100%).	82.14%
Ð	PASS	Verification	Completeness	None	FIELD	The number and percent of records with a NULL value in the visit_detail_id of the MEASUREMENT. (Threshold=100%).	80.90%
Ð	PASS	Verification	Completeness	None	FIELD	The number and percent of records with a NULL value in the value_source_value of the MEASUREMENT. (Threshold=100%).	79.89%
Ð	PASS	Validation	Completeness	None	TABLE	The number and percent of persons in the CDM that do not have at least one record in the DEVICE_EXPOSURE table (Threshold=100%).	76.70%
Ð	FAIL	Verification	Plausibility	Atemporal	CONCEPT	For the combination of CONCEPT_ID 3016049 (Testosterone Free [Mass/volume] in Serum or Plasma) and UNIT_CONCEPT_ID 8845 (picogram per milliliter), the number and percent of records that have value less than 5.00e+00. (Threshold=1%).	72.43%
Show	ing 126	to 130 of 3,35	1 entries			Previous 1 25 26 27	671 Next









IBM® MARKETSCAN® MULTI-STATE MEDICAID DATABASE

OVERVIEW

METADATA

RESULTS

ABOUT

DATA QUALITY ASSESSMENT

IBM® MARKETSCAN® MULTI-STATE MEDICAID DATABASE

Results generated at 2020-08-24 15:44:34 in 3 hours

		Ver	ification			Va	lidation				Total	
	Pass	Fail	Total	% Pass	Pass	Fail	Total	% Pass	Pass	Fail	Total	% Pass
Plausibility	1849	6	1855	100%	281	6	287	98%	2130	12	2142	99%
Conformance	550	13	563	98%	80	0	80	100%	630	13	643	98%
Completeness	322	5	327	98%	12	0	12	100%	334	5	339	99%
Total	2721	24	2745	99%	373	6	379	98%	3094	30	3124	99%

https://data.ohdsi.org/DataQualityDashboardMDCD/



> J Am Med Inform Assoc. 2021 Sep 18;28(10):2251-2257. doi: 10.1093/jamia/ocab132.

Increasing trust in real-world evidence through evaluation of observational data quality

Clair Blacketer ¹², Frank J Defalco ¹, Patrick B Ryan ¹³, Peter R Rijnbeek ²

Affiliations + expand PMID: 34313749 PMCID: PMC8449628 DOI: 10.1093/jamia/ocab132 Free PMC article

Abstract

Objective: Advances in standardization of observational healthcare data have enabled methodological breakthroughs, rapid global collaboration, and generation of real-world evidence to improve patient outcomes. Standardizations in data structure, such as use of common data models, need to be coupled with standardized approaches for data quality assessment. To ensure confidence in real-world evidence generated from the analysis of real-world data, one must first have confidence in the data itself.

Materials and methods: We describe the implementation of check types across a data quality framework of conformance, completeness, plausibility, with both verification and validation. We illustrate how data quality checks, paired with decision thresholds, can be configured to customize data quality reporting across a range of observational health data sources. We discuss how data



ACHILLES

Anthony Molinaro CDM Workshop 2022 2022-03-15



ETL Process





Data experts and CDM experts together design the ETL People with medical knowledge create the code mappings



ETL













Achilles is a data characterization and quality tool available for download here: https://github.com/OHDSI/Achilles

To visualize the results, a new tool was developed called ARES

https://github.com/OHDSI/Ares



2000 2001

date



PERSON REPORT





Population by Year of Birth



\bigcirc	Report Category		Data Source		Data Source Release		Report
77	Data Source Release	•	IBM CCAE	•	\$ 2022-01-22	•	Visit Occurrence

VISIT OCCURRENCE

Q Search in Table Ш Снооз	E COLUMNS TO DISPLAY
---------------------------	----------------------

Concept Id	Concept Name	↓ % People
9202	Outpatient Visit	80.64 %
581458	Pharmacy visit	63.08 %
9203	Emergency Room Visit	23.88 %
32036	Laboratory Visit	16.86 %
9201	Inpatient Visit	13.06 %

Repo Data	port Category ata Source Release					•		9	Data IBN	a Sou 1 CC	irce AE				•		0	Data 2022	Sourc 2-01-	ce Re •22	lease			•		1	Repo Con	ort ditio	ons	
5 <mark>1</mark> fier			N	2 umb	1,5 er of	<mark>85,88</mark> f Peo	9 ple				%	% 5 of ₽	<mark>1%</mark> Peopl	e				Reco	nds p	2.3 Der Pe	erson	Ē								
	2:-																													
Irst I	Jiag	gno	SIS	5																										
					+	+	_	-	_	-															+		_	_		
	51 fier	51 fier	51 fier irst Diagno	51 Provide Addition of the second sec	51 Sector of the	51	51	51 Provide a construction of the second	1 Image: Sector of the sec	1 Image: Control of the second se	1 1,585,889 fier Number of People	Image: Second constraint Image: Second const	Image: Sector of the sector of th	Image: Sector of the sector of th	on contraction of People % of People irst Diagnosis	Image: Sector of the sector	61 fier 1,585,889 Number of People % of People irst Diagnosis	51 A 1,585,889 % 1% fier Number of People % of People irst Diagnosis Image: Control of the second se	51 A Control 1,585,889 fier Number of People % of People Reco irst Diagnosis	Image: State of the state Image: State of the state Image: State of the state Image: State of the state Image: State of the state Image: State of the state Image: State of the state Image: State of the state Image: State Image: State of the state Image: State Image: State Image: State Image: State Image: State	51 A: 1,585,889 % 1% 2.3 fier Number of People % of People Records per P irst Diagnosis Image: State of the state	51 All 1,585,889 % 1% m 2.3 fier Number of People % of People Records per Person irst Diagnosis Image: Second sec	51 Provide Market 52 Provide Market 53 Provide Market 54 Provide Market 55 Provide Market 56 Provide Market 57 Provide Market 58 Provide Market 59 Provide Market 50 Provide Market 51 Provide Market 52 Provide Market 53 Provide Market 54 Provide Market 55 Provide Market 56 Provide Market 57 Provide Market 58 Provide Market 57 Provide Market 57 Provide Market 58 Provide Market 58 Provide Market 58 <td>Image: Second state Image: Second state Image: Second state Image: Second state<td>51 Image: Sector delta del</td><td>51 61 62 63 64 65 7 7 <td< td=""><td>51 All 1,585,889 X 1% Image: Contraction of the second seco</td><td>51 2.3 fier X 1% Number of People % of People rest Diagnosis</td><td>1 23 1 1595,889 Number of People % of People 23 Records per Person irst Diagnosis</td><td>Image: State of the state Image: State</td></td<></td></td>	Image: Second state Image: Second state Image: Second state Image: Second state <td>51 Image: Sector delta del</td> <td>51 61 62 63 64 65 7 7 <td< td=""><td>51 All 1,585,889 X 1% Image: Contraction of the second seco</td><td>51 2.3 fier X 1% Number of People % of People rest Diagnosis</td><td>1 23 1 1595,889 Number of People % of People 23 Records per Person irst Diagnosis</td><td>Image: State of the state Image: State</td></td<></td>	51 Image: Sector delta del	51 61 62 63 64 65 7 7 <td< td=""><td>51 All 1,585,889 X 1% Image: Contraction of the second seco</td><td>51 2.3 fier X 1% Number of People % of People rest Diagnosis</td><td>1 23 1 1595,889 Number of People % of People 23 Records per Person irst Diagnosis</td><td>Image: State of the state Image: State</td></td<>	51 All 1,585,889 X 1% Image: Contraction of the second seco	51 2.3 fier X 1% Number of People % of People rest Diagnosis	1 23 1 1595,889 Number of People % of People 23 Records per Person irst Diagnosis	Image: State of the state Image: State

- MEDIAN_VALUE: 39 P75_VALUE: 52 P90_VALUE: 59 MAX_VALUE: 65