



Rensselaer

why not change the world?®



Tetherless World Constellation

SciKG Part 3: Semantic Data Dictionaries (SDD)

Henrique Santos, Paulo Pinheiro, Jamie P. McCusker, Sabbir M. Rashid, Deborah L. McGuinness
May 28th 2023

The 20th Extended Semantic Web Conference (ESWC-23)

May 28th, 2023

SciKG Part 3: Semantic Data Dictionaries (SDDs)

Henrique Santos

Paulo Pinheiro

Jamie P. McCusker

Sabbir M. Rashid

Deborah L. McGuinness



Columbia University
MAILMAN SCHOOL
OF PUBLIC HEALTH

INTRODUCTION

We introduce this work,
motivate its importance, and
present our claims

01

RELATED WORK

We review literature related to
traditional data dictionaries, data
integration, mapping languages,
and semantic ETL

02

SEMANTIC DATA DICTIONARY

We present the various components
included in the Semantic Data
Dictionary specification

03

TABLE OF CONTENTS

04

MODELING APPROACHES

We discuss some modeling
strategies and provide some
examples to help illustrate this
work

05

CHALLENGES

We discuss some challenges
faced by domain scientists
when creating their of
Semantic Data Dictionaries

06

CONCLUSION

Thanks for listening!



01

INTRODUCTION

We introduce this work, motivate its importance, and present our claims

Data Int. journal paper authors - <https://bit.ly/3kG6iDi>

SABBIR RASHID



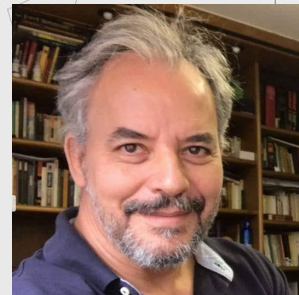
JAMIE MCCUSKER



PAULO PINHEIRO



MARCELLO BAX



HENRIQUE SANTOS



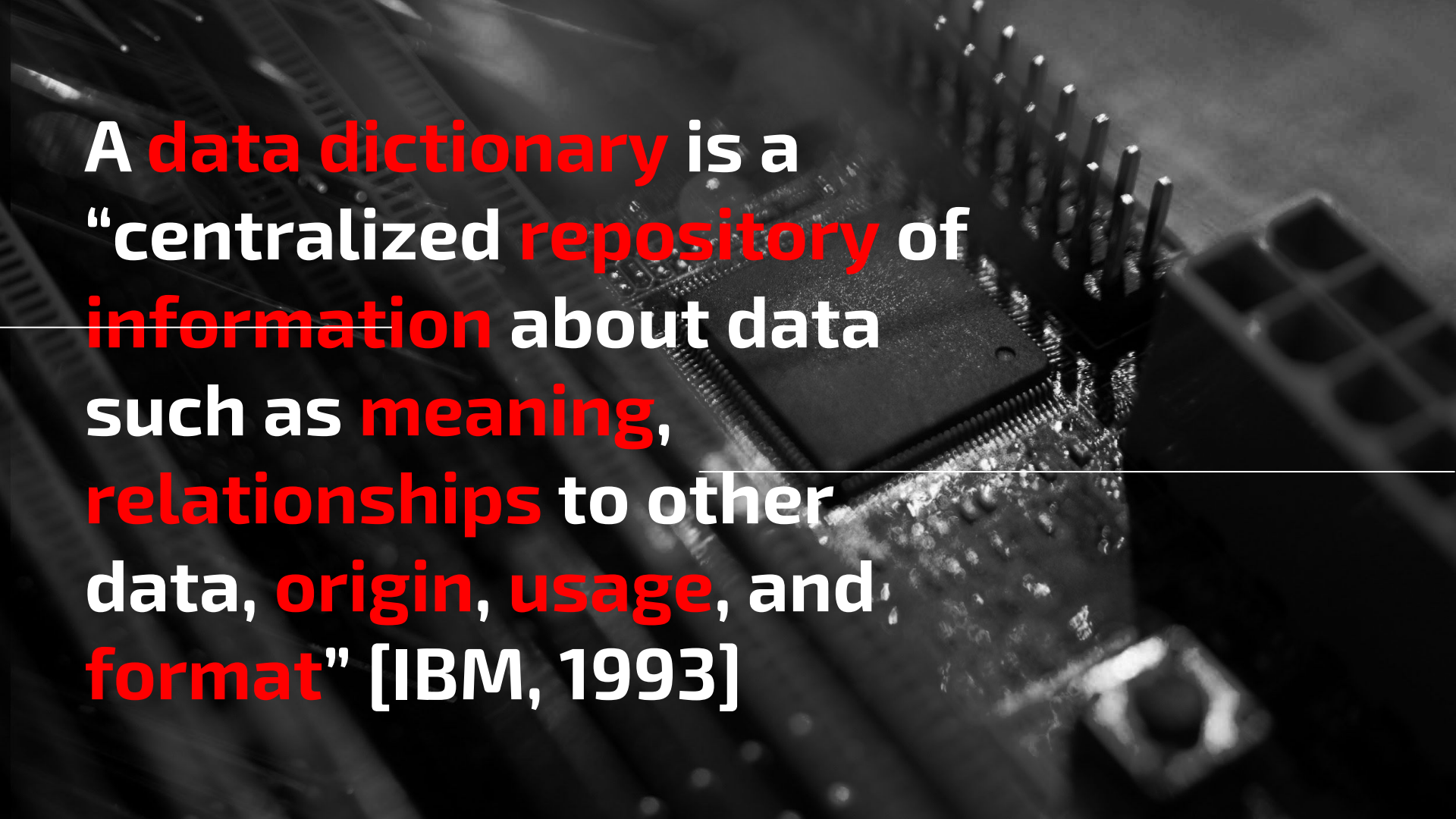
JEANETTE STINGONE



AMAR DAS



DEBORAH MCGUINNESS



A **data dictionary** is a
“centralized **repository** of
information about data
such as **meaning**,
relationships to other
data, **origin**, **usage**, and
format” [IBM, 1993]

DATA DICTIONARIES

- Ambiguity
- Standard adherence
- Human consumption

LIMITATIONS

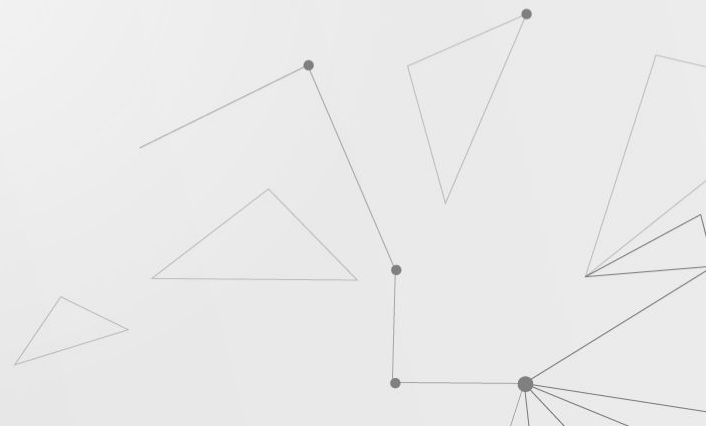
- Semantic technology usage
- Implicit concept annotation
- Provenance incorporation

SEMANTIC IMPROVEMENTS



MOTIVATION

- Annotate data from various domains
 - Harmonize data from multiple sources
 - Understand the data
-



CLAIMS

ADDRESSES LIMITATIONS

Addresses limitations
of traditional data
dictionaries

Semantic Data Dictionary

An approach for annotating and
transforming data

ABSTRACTION

Presents a level of
abstraction over mapping
language-based approaches

F.A.I.R.

Resulting model is
Findable, Accessible,
Interoperable, and
Reusable





02

RELATED WORK

We review literature related to traditional data dictionaries, data integration, mapping languages, and semantic ETL

TRADITIONAL DATA DICTIONARIES

LIMITATIONS

- Data dictionaries mentioned in patents [Haskell et al., 2009, Lau et al., 2002, Apacible et al., 2013]
 - Stony Brook Data Governance Council (<https://bit.ly/3oD4g90>)
 - The Open Science Framework (<https://bit.ly/35EPupT>)
 - Biosystematic Database of World Diptera [Thompson, 1999]
 - Project Open Data Metadata Schema (<https://bit.ly/3oCYTqr>)
- Minimal incorporation of semantics
 - Object and relation elicitation not permitted
 - Domain-specific
 - Not machine-readable
 - Lack of a formal creation standard

data/

OSF Storage (United States)

0.1 datadictionary.csv

0.2 raw.pdf

0.3 clean.csv

Sheet_1

Show rows with cells including:

Variable	Variable name	Mesaurement unit	Allowed values	Description
Participant ID number	ID	Numeric	001-999	ID number assigned to participant in sequential order
Group number	GROUP	Numeric	1-30	Group assigned to participant based on ID number
Age in years	AGE	Numeric	18.0-65.0	Age of participant in years
Date of birth	DOB	mm/dd/yyyy	1-12/1-31/1951-1998	Participant's date of birth
Gender	SEX	Numeric	1 = male 2 = female	Participant's gender
Date of survey	SURVEY	mm/dd/yyyy	01/01/2015 – 01/01/2016	When the participant completed the survey
Self-reported consumer spending	SPEND	Numeric	0-100,000,000	Self-reported average yearly expenditure
Market sentiment	SENTIMENT	Numeric	1 = negative 2 = neutral 3 = positive	Sentiment towards US domestic economy
Actual GDP growth	GDP	Numeric	-5.0-5.0	Average US yearly GDP growth

Tags

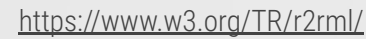
DATA INTEGRATION



- Techniques that utilize data from multiple sources to construct a unified view of the combined data [Lenzerini, 2002]
- The Semantic Web Integration Tool (SWIT) [del Carmen Legaz-Garcia et al., 2016]
- RDF-Gen [Santipantakis et al., 2018]
- DataOps [Pinkel et al., 2015]
- OpenRefine [Ham, 2013]

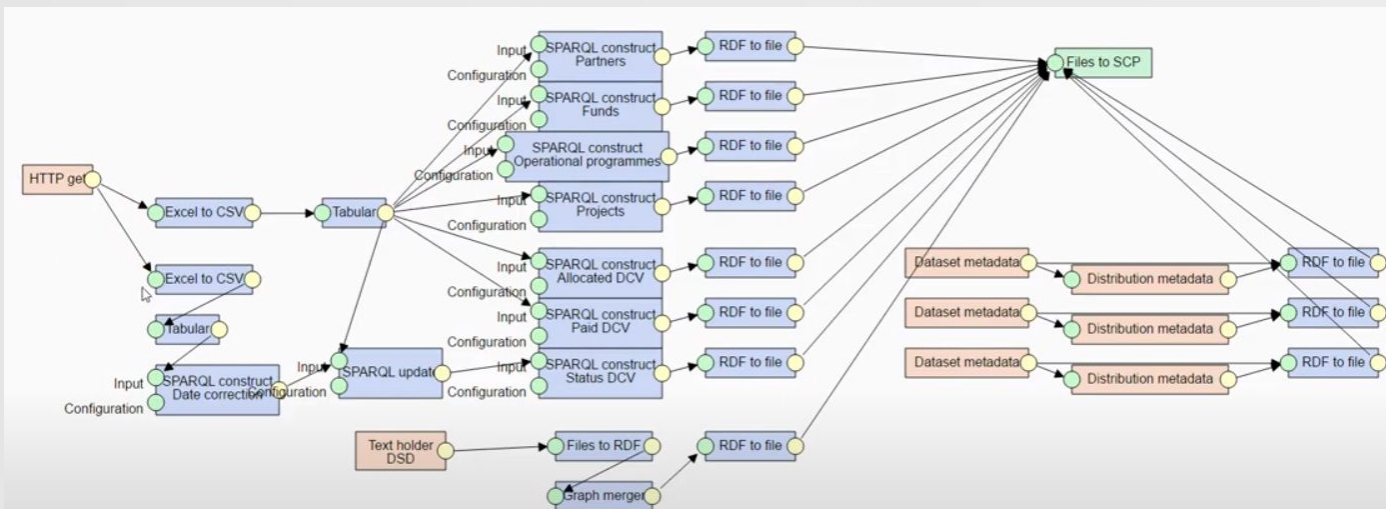
LIMITATIONS

- Not all tools are open source
- Some require knowledge of mapping languages
- Difficulties with subset selection, cell-based operations, dataset merging
- Not all tools allow object elicitation
- Some difficulties associated with adoption



<https://www.w3.org/TR/r2rml/>

- <https://etl.linkedpipes.com/>





03

SEMANTIC DATA DICTIONARY

We present the various components
included in the Semantic Data
Dictionary specification

SEMANTIC DATA DICTIONARY SPECIFICATION



INFOSHEET

Contains links to the other specifications



DICTIONARY MAPPING

Used to annotate the columns of a dataset



CODEBOOK

Used to annotate coded values



METADATA SUPPLEMENT

Includes metadata about the Semantic Data Dictionary or the associated dataset



CODE MAPPING

Used to encode shortcut notations
(See <https://bit.ly/2HLydmK>)



TIMELINE

Used for complex temporal mappings



PROPERTIES

Used to customize the properties used during the mapping process



INFOSHEET

SPECIFICATION

METADATA SUPPLEMENT

Infosheet Row	Description
CODE MAPPING	Reference to Code Mapping table location
CODEBOOK	Reference to Codebook table location
DICTIONARY MAPPING	Reference to Dictionary Mapping table location
PROPERTIES	Reference to Properties table location
TIMELINE	Reference to Timeline table location

Infosheet Row	Related Property	Description
CONTRIBUTORS	<i>dct:contributor</i>	Contributors to the SDD
CREATORS	<i>dct:creator</i>	Creators of the SDD
DATE CREATED	<i>dct:created</i>	Date the SDD was created
DESCRIPTION	<i>dct:description</i>	Description of the KG fragment
IMPORTS	<i>owl:imports</i>	Ontologies that the SDD references
KEYWORDS	<i>schema:keywords</i>	Keywords to be associated with the KG fragment
LICENSE	<i>dct:license</i>	License URL
PREVIOUS VERSION	<i>pav:previousVersion</i>	Previous version URL
PUBLISHER	<i>dct:publisher</i>	Publisher of the SDD
TITLE	<i>dct:title</i>	Title of KG fragment
VERSION	<i>owl:versionInfo</i>	Current version URL
VERSION OF	<i>dct:isVersionOf</i>	Resource URL for primary version

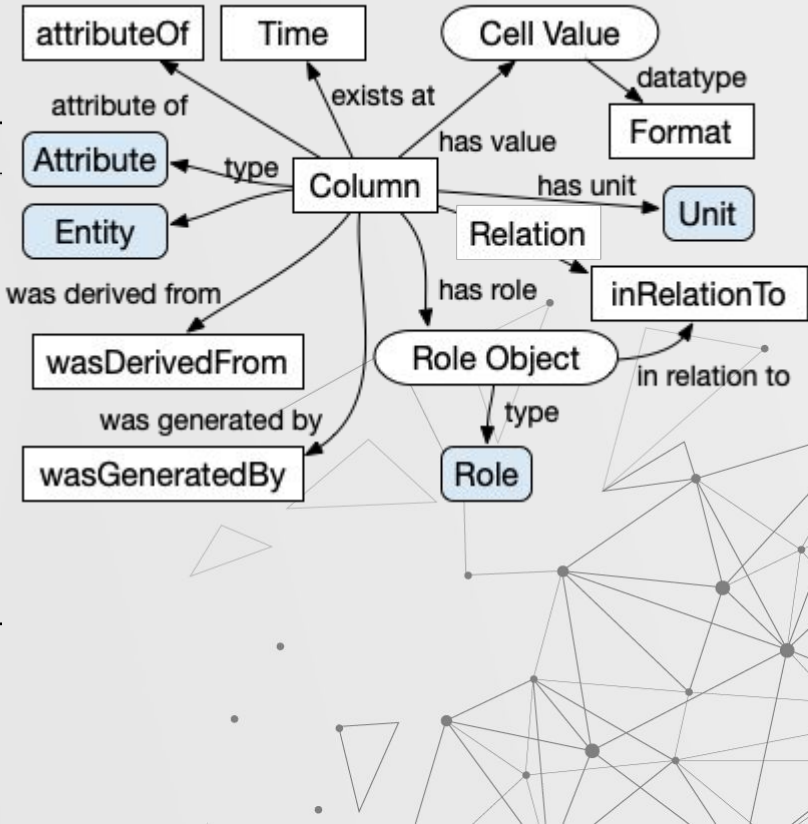
- Data on the Web Best Practices
 - <https://www.w3.org/TR/dwbp/>
- Semantic Web Health Care and Life Sciences
 - <https://www.w3.org/TR/hcls-dataset/>

DICTIONARY MAPPING

DIAGRAM

SPECIFICATION

DM Column	Related Property	Description
ATTRIBUTE	<i>rdf:type</i>	Class of attribute entry
ATTRIBUTEOF	<i>sio:isAttributeOf</i>	Entity having the attribute
COLUMN		Entry column header in dataset
ENTITY	<i>rdf:type</i>	Class of entity entry
FORMAT		Specifies the structure of the cell value
INRELATIONTO	<i>sio:inRelationTo</i>	Entity that the role is linked to
LABEL	<i>rdfs:label</i>	Label for the entry
RELATION		Custom property used in INRELATIONTO
ROLE	<i>sio:hasRole</i>	Type of the role of the entry
TIME	<i>sio:existsAt</i>	Time point of measurement
UNIT	<i>sio:hasUnit</i>	Unit of measure for entry
WASDERIVEDFROM	<i>prov:wasDerivedFrom</i>	Entity from which the entry was derived
WASGENERATEDBY	<i>prov:wasGeneratedBy</i>	Activity from which the entry was produced



DICTIONARY MAPPING FORMALISM

$\exists \text{COLUMN} \wedge \exists \text{ATTRIBUTE} \Rightarrow \text{ATTRIBUTE}(\text{COLUMN})$

$\exists \text{COLUMN} \wedge \exists \text{ENTITY} \Rightarrow \text{ENTITY}(\text{COLUMN})$

$\exists \text{COLUMN} \wedge \exists \text{LABEL} \Rightarrow \text{rdfs:label}(\text{COLUMN}, \text{LABEL})$

$\exists \text{COLUMN} \wedge \exists \text{COMMENT} \Rightarrow \text{rdfs:comment}(\text{COLUMN}, \text{COMMENT})$

$\exists \text{COLUMN} \wedge \exists \text{DEFINITION} \Rightarrow \text{skos:definition}(\text{COLUMN}, \text{DEFINITION})$

$\exists \text{COLUMN} \wedge \exists \text{ATTRIBUTEOF} \Rightarrow \text{sio:attributeOf}(\text{COLUMN}, \text{ATTRIBUTEOF})$

$\exists \text{COLUMN} \wedge \exists \text{UNIT} \Rightarrow \exists U \wedge \text{UNIT}(U) \wedge \text{sio:hasUnit}(\text{COLUMN}, U)$

$\exists \text{COLUMN} \wedge \exists \text{FORMAT} \wedge \exists \text{Value} \Rightarrow \text{sio:hasValue}(\text{COLUMN}, \text{Value}^{\sim \text{FORMAT}})$

$\exists \text{COLUMN} \wedge \exists \text{TIME} \Rightarrow \text{sio:existsAt}(\text{COLUMN}, \text{TIME})$

$\exists \text{COLUMN} \wedge \exists \text{ROLE} \Rightarrow \exists R \wedge \text{sio:hasRole}(\text{COLUMN}, R) \wedge \text{ROLE}(R)$

$\exists \text{COLUMN} \wedge \exists \text{ROLE} \wedge \exists \text{INRELATIONTO} \Rightarrow \exists R \wedge \text{sio:hasRole}(\text{COLUMN}, R) \wedge \text{ROLE}(R) \\ \wedge \text{sio:inRelationTo}(R, \text{INRELATIONTO})$

$\exists \text{COLUMN} \wedge \exists \text{INRELATIONTO} \Rightarrow \text{sio:inRelationTo}(\text{COLUMN}, \text{INRELATIONTO})$

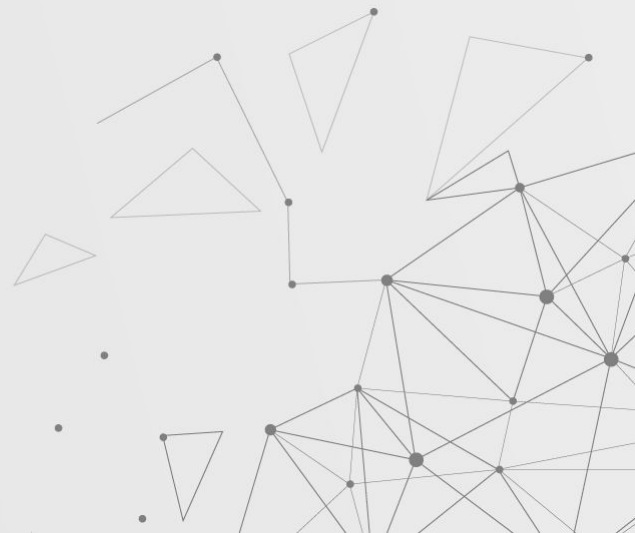
$\exists \text{COLUMN} \wedge \exists \text{RELATION} \wedge \exists \text{INRELATIONTO} \Rightarrow \text{RELATION}(\text{COLUMN}, \text{INRELATIONTO})$

$\exists \text{COLUMN} \wedge \exists \text{ROLE} \wedge \exists \text{RELATION} \wedge \exists \text{INRELATIONTO} \Rightarrow \exists R \wedge \text{sio:hasRole}(\text{COLUMN}, R) \wedge \text{ROLE}(R) \\ \wedge \text{RELATION}(R, \text{INRELATIONTO})$

$\exists \text{COLUMN} \wedge \exists \text{WASDERIVEDFROM} \Rightarrow \text{prov:wasDerivedFrom}(\text{COLUMN}, \text{WASDERIVEDFROM})$

$\exists \text{COLUMN} \wedge \exists \text{WASGENERATEDBY} \Rightarrow \text{prov:wasGeneratedBy}(\text{COLUMN}, \text{WASGENERATEDBY})$

$\exists \text{COLUMN} \wedge \exists \text{Value} \Rightarrow \text{sio:hasValue}(\text{COLUMN}, \text{Value})$

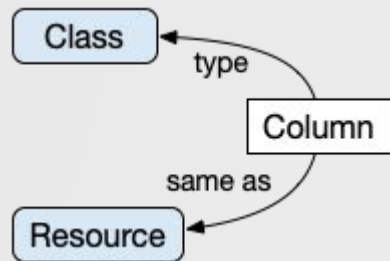


CODEBOOK

SPECIFICATION

Codebook Column	Related Property	Description
CLASS	<i>rdf:type</i>	Class the Code refers to
CODE	<i>sio:hasValue</i>	Value of the dataset entry
COLUMN		Entry column header in dataset
LABEL	<i>rdfs:label</i>	Label for the codebook entry
RESOURCE	<i>rdf:type</i>	Web Resource URI the Code refers to

DIAGRAM



FORMALISM

$\exists \text{COLUMN} \wedge \exists \text{CLASS} \Rightarrow \text{CLASS}(\text{COLUMN})$

$\exists \text{COLUMN} \wedge \exists \text{LABEL} \Rightarrow \text{rdfs:label}(\text{COLUMN}, \text{LABEL})$

$\exists \text{COLUMN} \wedge \exists \text{RESOURCE} \Rightarrow \text{owl:sameAs}(\text{COLUMN}, \text{RESOURCE})$

$\exists \text{COLUMN} \wedge \exists \text{CODE} \Rightarrow \text{sio:hasValue}(\text{COLUMN}, \text{CODE})$

TIMELINE

SPECIFICATION

Timeline Column	Related Property	Description
END	<i>sio:hasEndTime</i>	The starting time point associated with the Timeline entry
INRELATIONTO	<i>sio:inRelationTo</i>	Entity that the Timeline entry is associated with
NAME		Implicit entry reference for the Timeline entry
START	<i>sio:hasStartTime</i>	The starting time point associated with the Timeline entry
TYPE	<i>rdf:type</i>	Class the Timeline entry refers to
UNIT	<i>sio:hasUnit</i>	Unit of measure for Timeline entry

FORMALISM

$$\exists \text{NAME} \wedge \exists \text{TYPE} \Rightarrow \text{TYPE}(\text{NAME})$$

$$\exists \text{NAME} \wedge \exists \text{LABEL} \Rightarrow \text{rdfs:label}(\text{NAME})$$

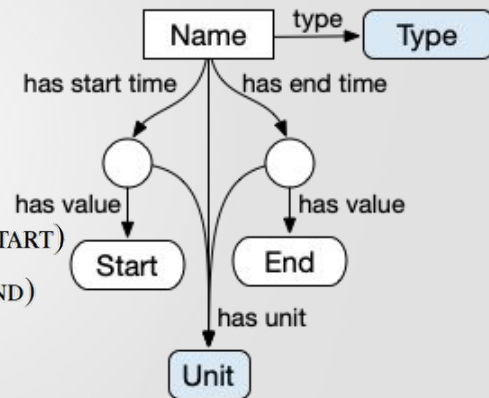
$$\exists \text{NAME} \wedge \exists \text{START} \Rightarrow \exists \text{S} \wedge \text{sio:hasStartTime}(\text{NAME}, \text{S}) \wedge \text{sio:hasValue}(\text{S}, \text{START})$$

$$\exists \text{NAME} \wedge \exists \text{END} \Rightarrow \exists \text{E} \wedge \text{sio:hasEndTime}(\text{NAME}, \text{E}) \wedge \text{sio:hasValue}(\text{E}, \text{END})$$

$$\exists \text{NAME} \wedge \exists \text{START} \wedge \exists \text{END} \wedge \text{START} \equiv \text{END} \Rightarrow \exists \text{T} \wedge \text{sio:existsAt}(\text{NAME}, \text{T}) \wedge \text{sio:hasValue}(\text{T}, \text{START})$$

$$\exists \text{NAME} \wedge \exists \text{UNIT} \Rightarrow \exists \text{U} \wedge \text{UNIT}(\text{U}) \wedge \text{sio:hasUnit}(\text{NAME}, \text{U})$$

$$\exists \text{NAME} \wedge \exists \text{INRELATIONTO} \Rightarrow \text{sio:inRelationTo}(\text{NAME}, \text{INRELATIONTO})$$



DIAGRAM

PROPERTIES SPECIFICATION

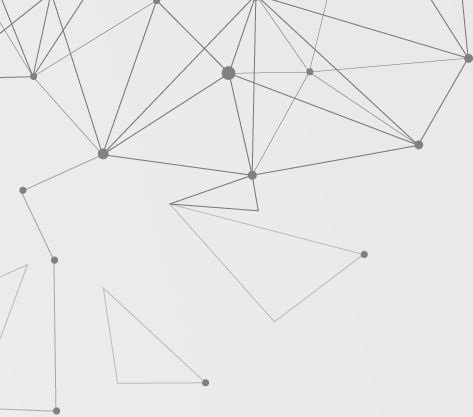
Row	Property
ATTRIBUTE	<i>rdf:type</i>
ATTRIBUTEOf	<i>sio:isAttributeOf</i>
COMMENT	<i>rdfs:comment</i>
DEFINITION	<i>skos:definition</i>
END	<i>sio:hasEndTime</i>
ENTITY	<i>rdf:type</i>
INRELATIONTo	<i>sio:inRelationTo</i>
LABEL	<i>rdfs:label</i>
ROLE	<i>sio:hasRole</i>
START	<i>sio:hasStartTime</i>
TIME	<i>sio:existsAt</i>
TYPE	<i>rdf:type</i>
UNIT	<i>sio:hasUnit</i>
VALUE	<i>sio:hasValue</i>
WASDERIVEDFROM	<i>prov:wasDerivedFrom</i>
WASGENERATEDBy	<i>prov:wasGeneratedBy</i>



04

MODELLING APPROACHES

We discuss some modeling strategies
and provide some examples to help
illustrate this work



SEMANTIC SCIENTIFIC WORKFLOW





ONTOLOGY ENGINEERING

CLASS SELECTION

- Collect relevant ontologies
 - <http://www.ontobee.org/>
 - <https://bioportal.bioontology.org/>


SUPPORTING ONTOLOGY

- What if concepts used to annotate dataset does not exist in an ontology?
 - Create concept map
 - Engineer a supporting ontology
 - Protege
 - <https://protege.stanford.edu/>
 - Manual vs. automated approaches





KNOWLEDGE GRAPH GENERATION

- <https://github.com/tetherless-world/SemanticDataDictionary>
 - `sdd2rdf`
 - `sdd2setl`
 - <https://github.com/tetherless-world/whyis>
 - Will cover this in the tutorial!
 - Loading of KG into a triplestore
 - Querying the resulting Graph
- 

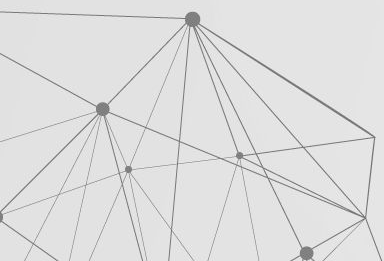
INFOSHEET EXAMPLE

Attribute	Value
CREATORS	Sabbir M. Rashid
CODE MAPPING	NHANES/config/code_mappings.csv
CODEBOOK	NHANES/input/CB/DEMO_H.Doc-CB.csv
CONTRIBUTORS	“James P. McCusker, Paulo Pinheiro, Marcello P. Bax, Henrique O. Santos, Alexander New, Shruthi Chari, Mathew Johnson, John S. Erickson, Kristin P. Bennett, Jeanette A. Stingone, Deborah L. McGuinness”
DATE CREATED	2018-10-14
DESCRIPTION	KG fragment from manually annotated NHANES Demographics SDD.
DICTIONARY MAPPING	NHANES/input/DM/DEMO_H.Doc-DM.csv
IMPORTS	“ http://semanticscience.org/ontology/sio-subset-labels.owl , http://hadatac.org/ont/chear/ , http://purl.obolibrary.org/obo/ncit.owl ”
KEYWORDS	“demographics, gender, age, race, citizenship, marital status, household”
LICENSE	https://opensource.org/licenses/MIT
PREVIOUS VERSION	http://tw.rpi.edu/heals/kb/nhanes/1.1
PROPERTIES	NHANES/config/Properties.csv
PUBLISHER	Tetherless World Constellation
TIMELINE	NHANES/input/TL/DEMO_H.Doc-TL.csv
TITLE	The National Health and Nutrition Examination (NHANES) SDD KG
VERSION	http://tw.rpi.edu/heals/kb/nhanes/1.2
VERSION OF	http://tw.rpi.edu/heals/kb/nhanes/

DICTIONARY MAPPING EXAMPLE

COLUMN	LABEL	ATTRIBUTE	ATTRIBUTEOF	UNIT	TIME	ENTITY	RELATION	INRELATIONTO
SEQN	Respondent sequence number	sio:Identifier	??participant					
RIAGENDR	Gender	sio:BiologicalSex	??participant					
RIDAGEYR	Age in years at screening	sio:Age	??participant	yr	??screening			
RIDAGEMN	Age in months at screening	sio:Age	??participant	moth	??screening			
RIDRETH1	Race/Hispanic origin	sio:Race	??participant					
RIDEXAGM	Age in months at exam	sio:Age	??participant	moth	??exam			
DMDBORN4	Country of birth				??birth	sio:Country	sio:isLocationOf	??participant
DMDCITZN	Citizenship status	sio:StatusDescriptor	??participant					
DMDYRSUS	Length of time in US	sio:TimeInterval	??participant					
DMDDEDUC3	Education level - Children/Youth	chear:EducationLevel	??participant					
DMDDEDUC2	Education level - Adults 20+	chear:EducationLevel	??participant					
DMDMAR1L	Marital status	chear:MaritalStatus	??participant					
RIDEXPRG	Pregnancy status at exam	sio:StatusDescriptor	??pregnancy		??exam			??participant
SIALANG	Language of SP Interview	chear:Language	??instrument		??interview			??participant
DMDHRGND	HH ref person's gender	sio:BiologicalSex	??HHRef					
DMDHRAGE	HH ref person's age in years	sio:Age	??HHRef	yr				
DMDHRBR4	HH ref person's country of birth				??birth	sio:Country	sio:isLocationOf	??HHRef
DMDHREDU	HH ref person's education level	chear:EducationLevel	??HHRef					
DMDHRMAR	HH ref person's marital status	chear:MaritalStatus	??HHRef					
WTINT2YR	Full sample 2 year interview wt	chear:Weight	??participant		??interview			
WTMEC2YR	Full sample 2 year MEC exam wt	chear:Weight	??participant		??exam			
INDHHIN2	Annual household income	chear:Income	??household					

EXPLICIT ENTRIES

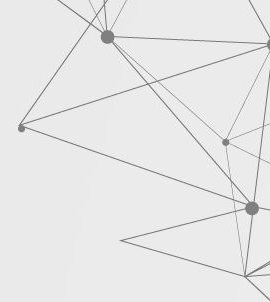


IMPLICIT ENTRIES

COLUMN	LABEL	ENTITY	ROLE	INRELATIONTO
??participant	Participant	ncit:C29867, sio:Human	sio:SubjectRole	
??screening	Screening	chear:Screening		
??exam	Examination	ncit:C131902		
??birth	Birth	sio:Birthing		
??pregnancy	Pregnancy	chear:Pregnancy		
??interview	Interview	ncit:C16751		
??instrument	Instrumentation	ncit:C16742		
??household	Household	chear:Household		??participant
??HHRef	Household reference	sio:Human	chear:HeadOfHousehold	??household

CODEBOOK EXAMPLE

COLUMN	CODE	LABEL	CLASS
RIAGENDR	1	Male	sio:Male
RIAGENDR	2	Female	sio:Female
RIAGENDR	.	Missing	ncit:C142610
RIDRETH1	1	Mexican American	exo:0000151
RIDRETH1	2	Other Hispanic	exo:0000145
RIDRETH1	3	Non-Hispanic White	exo:0000158
RIDRETH1	4	Non-Hispanic Black	exo:0000132
RIDRETH1	5	Other Race - Including Multi-Racial	exo:0000153
RIDRETH1	.	Missing	ncit:C142610
DMDEDUC3	0	Never attended / kindergarten only	cheat:NoFormalEducation
DMDEDUC3	1	1st grade	cheat:EducationGrade
DMDEDUC3	2	2nd grade	cheat:EducationGrade
DMDEDUC3	3	3rd grade	cheat:EducationGrade
DMDEDUC3	4	4th grade	cheat:EducationGrade
DMDEDUC3	5	5th grade	cheat:EducationGrade
DMDEDUC3	6	6th grade	cheat:EducationGrade
DMDEDUC3	7	7th grade	cheat:EducationGrade
DMDEDUC3	8	8th grade	cheat:EducationGrade



The background features a complex network of dark grey nodes connected by thin, light grey lines, creating a web-like structure. The nodes are of varying sizes and are distributed across the slide, with a higher concentration on the right side. The overall aesthetic is modern and technical.


05

CHALLENGES

We discuss some challenges faced by domain scientists when creating their of Semantic Data Dictionaries




EXPERIMENTAL SETUP

- Domain scientists were presented with initial training
 - Epidemiologists and biostatisticians
 - Supporting materials were developed in collaboration with a domain expert
 - Were made available to provide guidance and examples
 - A template for completing the Semantic Data Dictionary was provided
 - Included pre-populated fields for common demographic concepts
 - Such as age, race, and gender
 - A help document was created that included instructions and representations of more complex concepts
 - Measurements of environmental samples
 - Measurements of biological samples
 - Measurements taken at specific time-points
 - A practical workshop was held
 - A semantic scientist provided training in semantic representation to the domain scientists
 - Domain scientists completed at least one Semantic Data Dictionary for an epidemiologic study
- 



CHALLENGES

- Domain scientists had representation difficulties
 - Complex ideas (e.g. fasting blood glucose levels)
 - Implicit concepts
 - Uncommon representation in the public health domain
 - Not necessarily intuitive
 - Time associations
 - Determining best ontology term for annotation was not always clear
 - What if a term was not found in a supporting ontology?
 - Best way to represent concept in a semantically appropriate way
 - What other ontologies should be used?
 - Requirement for user to have some domain & ontology knowledge
 - Currently only supports annotation of tabular data
 - Annotation process is mostly manual
 - Documentation and tutorials can be improved
- 



06

CONCLUSIONS

Thanks for listening!



CLOSING COMMENTS

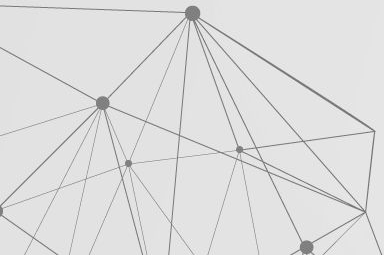
Semantic Data Dictionaries address many of the limitations of the prior work -- thus, this work helps advance the state-of-the-art

The SDD approach follows Semantic Web standards and results in artifacts that are findable, accessible, interoperable, and reusable



RESOURCES

- whyis - <https://github.com/tetherless-world/whyis>
- HADatAc - <https://github.com/paulopinheiro1234/hadatac>
- sdd2rdf - <https://github.com/tetherless-world/SemanticDataDictionary>
- Documentation - <https://tetherless-world.github.io/sdd/>
- Annotated resource examples - https://github.com/tetherless-world/sdd/tree/master/sdd_resources
- Journal Paper - <https://bit.ly/3kG6iDi>



REFERENCES

- Apacible, J. T., Nolan, S. P., Kalmady, G. D., and Varadan, V. (2013). Extensible and localizable health-related dictionary. US Patent 8,417,537.
- Arenas, M., Bertails, A., Prudhommeaux, E., and Sequeda, J. (2012). A direct mapping of relational data to rdf. W3C recommendation, 27:1–11.
- Cerans, K. and Bumans, G. (2011). Rdb2owl: a rdb-to-rdf/owl mapping specification language. Information Systems, pages 139–152.
- del Carmen Legaz-García, M., Minarro-Gimenez, J. A., Menarguez-Tortosa, M., and Fernández-Breis, J. T. (2016). Generation of open biomedical datasets through ontology-driven transformation and integration processes. Journal of biomedical semantics, 7(1):32.
- Dimou, A., Vander Sande, M., Colpaert, P., De Vocht, L., Verborgh, R., Mannens, E., and Van de Walle, R. (2014). Extraction and semantic annotation of workshop proceedings in html using rml. In Semantic Web Evaluation Challenge, pages 114–119. Springer.
- Ham, K. (2013). Openrefine (version 2.5). <http://openrefine.org>. free, open-source tool for cleaning and transforming data. Journal of the Medical Library Association: JMLA, 101(3): 233.
- Haskell, R. E., Heil, J. A., and Cassidy, J. (2009). Dynamic dictionary and term repository system. US Patent 7,580,831.
- IBM (1993). IBM Dictionary of Computing. McGraw-Hill, Inc., New York, NY, USA, 10th edition.
- Klimek, J., Skoda, P., and Necasky, M. (2016). Linkedpipes etl: Evolved linked data preparation. In European Semantic Web Conference, pages 95–100. Springer.
- Knoblock, C. A. and Szekely, P. (2015). Exploiting semantics for big data integration. AI Magazine, 36(1).
- Kuchinke, W., Wiegmann, S., Verplancke, P., and Ohmann, C. (2006). Extended cooperation in clinical studies through exchange of cdisc metadata between different study software solutions. Methods of information in medicine, 45(04):441–446.

REFERENCES

- Lau, L., Endo, J., Karren, S., Willis, M., Harada, S., Beeney, S., Larsen, B., Cassin, E., and Gerard, M. (2002). Mapping clinical data with a health data dictionary. US Patent App. 09/755,966.
- Lenzerini, M. (2002). Data integration: A theoretical perspective. In Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems, pages 233–246. ACM.
- McCusker, J. P., Chastain, K., Rashid, S., Norris, S., and McGuinness, D. L. (2018). Setlr: the semantic extract, transform, and load-r. PeerJ Preprints, 6:e26476v1.
- Michel, F., Djimenou, L., Zucker, C. F., and Montagnat, J. (2015). Translation of relational and non-relational databases into rdf with xr2rml. In 11th International Conference on Web Information Systems and Technologies (WEBIST'15), pages 443–454.
- Pinkel, C., Schwarte, A., Trame, J., Nikolov, A., Bastinos, A. S., and Zeuch, T. (2015). Dataops: seamless end-to-end anything-to-rdf data integration. In European Semantic Web Conference, pages 123–127. Springer.
- Post, A. R., Krc, T., Rathod, H., Agravat, S., Mansour, M., Torian, W., and Saltz, J. H. (2013). Semantic etl into i2b2 with eureka! AMIA Summits on Translational Science Proceedings, 2013:203.
- Santipantakis, G. M., Kotis, K. I., Vouros, G. A., and Doulkeridis, C. (2018). Rdf-gen: Generating rdf from streaming and archival data. In Proceedings of the 8th International Conference on Web Intelligence, Mining and Semantics, page 28. ACM.
- Schultz, A., Matteini, A., Isele, R., Bizer, C., and Becker, C. (2011). Ldif-linked data integration framework. In Proceedings of the Second International Conference on Consuming Linked Data-Volume 782, pages 125–130. CEUR-WS. Org.
- Slepicka, J., Yin, C., Szekely, P. A., and Knoblock, C. A. (2015). Kr2rml: An alternative interpretation of r2rml for heterogeneous sources. In Cold.

REFERENCES

- Stadler, C., Unbehauen, J., Westphal, P., Sherif, M. A., and Lehmann, J. (2015). Simplified rdb2rdf mapping. In LDOW@ WWW, volume 1409.
- Thompson, F. C. (1999). Data dictionary and standards for fruit fly information database.
- Myia. Warzel, D. B., Andonyadis, C., McCurry, B., Chilukuri, R., Ishmukhamedov, S., and Covitz, P. (2003). Common data element (cde) management and deployment in clinical trials. In AMIA annual symposium proceedings, volume 2003, page 1048. American Medical Informatics Association.
- Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.-W., da Silva Santos, L. B., Bourne, P. E., et al. (2016). The fair guiding principles for scientific data management and stewardship. Scientific data, 3.
- Zozus, M. N., Bonner, J., and Rock, L. (2017). Towards data value-level metadata for clinical studies. In ITCH, pages 418–423.