



# Aligning Diverse Data for a Personalized Health Knowledge Graph

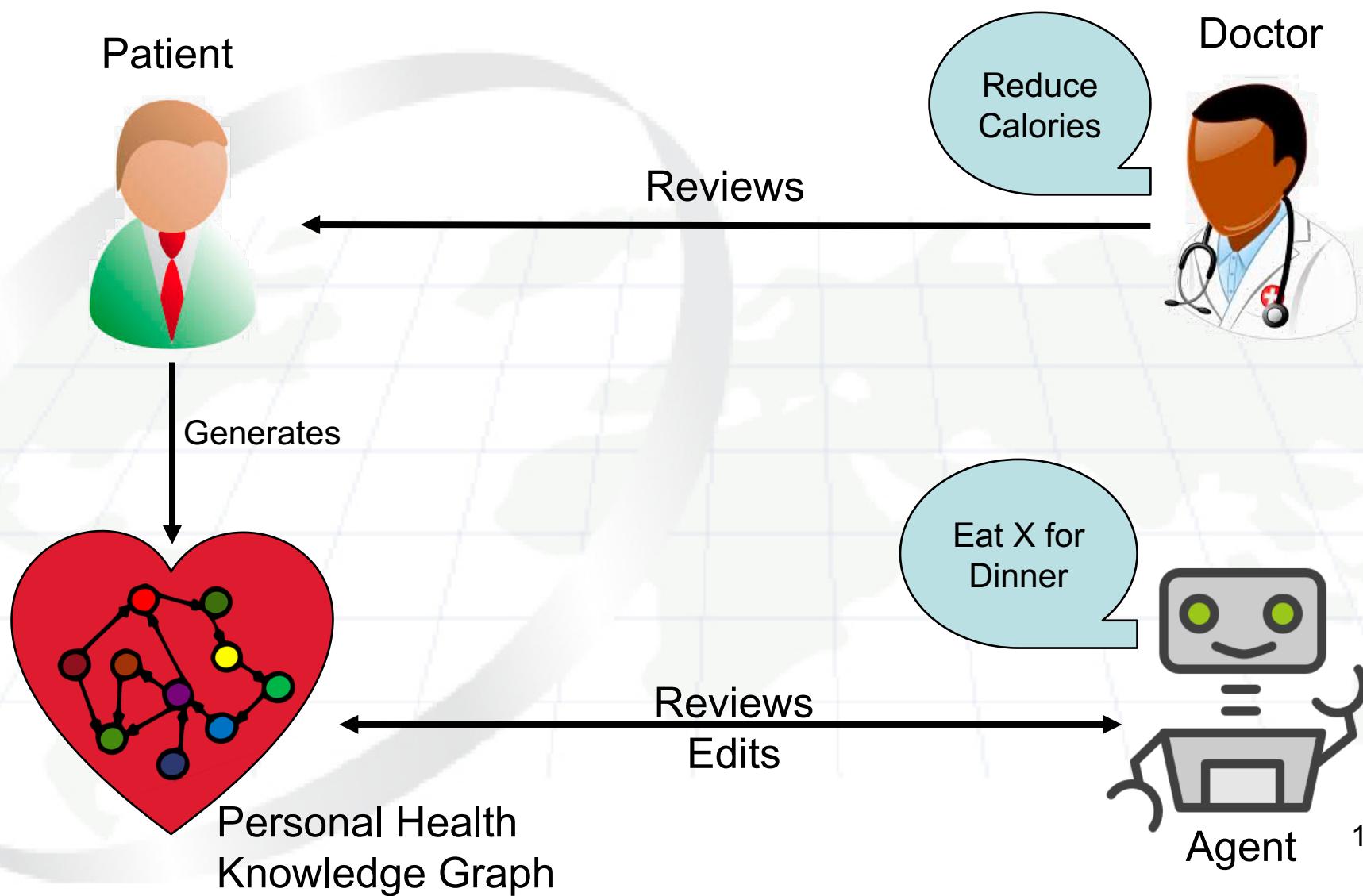
Matthew Johnson, Meenu Ravi, Sabbir Rashid,  
Paulo Pinheiro, and Deborah L. McGuinness

Rensselaer Polytechnic Institute, Troy, NY, 12180, USA



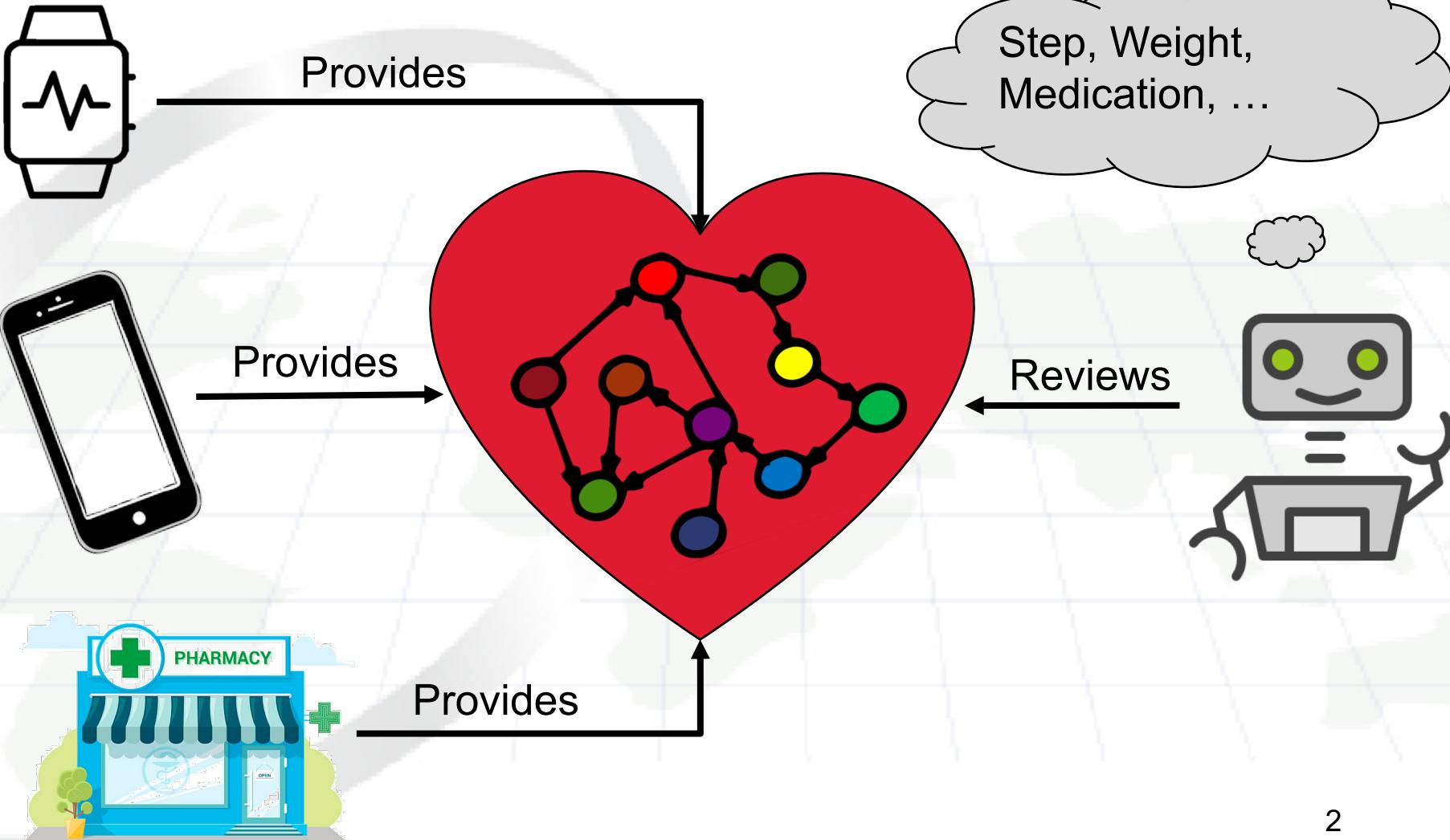


# Personal Health Knowledge Graph





# Integration Problem





# Semantic Data Dictionary

Data

id	age	gender	LNAS_C1
3479	23	...	...
...			

Data Dictionary

<b>id</b>	Participant ID number
<b>age</b>	Age of mother at enrollment
<b>gender</b>	Child gender
<b>LNAS_C1</b>	Natural log of drinking water arsenic at enrollment

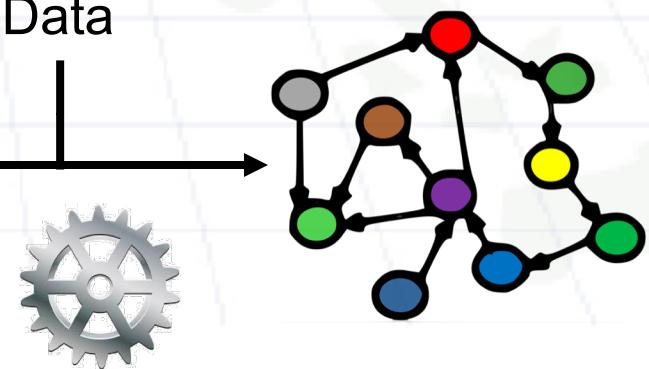


Column	Attribute	attributeOf	Unit
id	hasco:originalID	??mother	
age	sio:SIO_001013	??mother	sio:SIO_000428
gender	sio:SIO_010029	??child	

Virtual Column	Entity	Role	inRelationTo
??mother	sio:SIO_000485	hhear:00502	??child
??child	sio:SIO_000485	hhear:00492	??mother

Data

Knowledge Graph





# SDD-Editor

## Data Dictionary Column Description

SDD-34.xlsx

Save Download Populate Undo Redo Show Labels Shop Terms View Cart Verify SDD External Verify

Cell Value: sio:SIO\_001013

Column Description: Age of mother at enrollment

Column	Attribute	attributeOf	Unit	Time	Entity	Role	Relation	inRelationTo
1	id	hasco:originalID	??mother					??study
2	age	sio:SIO_001013	??mother	sio:SIO_000428	??pregnancy			
3	gender	sio:SIO_001014	??child					
4	??mother			sio:SIO_000485 +	hhear:00502			??child
5	??child			sio:SIO_000485	hhear:00492			??mother

Suggestion Service: ✓ Refresh Suggestion: C

Choose from Below

- 1.00 age (sio)
- 0.85 age (obo)
- 0.75 Child (hhear)
- 0.73 Birth (hhear)

Results Found

Age

-Description: How long something has existed; elapsed time since birth.  
<http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#C25150>

Age

-Description: How long something has existed; elapsed time since birth.  
[http://purl.obolibrary.org/obo/NCIT\\_C25150](http://purl.obolibrary.org/obo/NCIT_C25150)

Age

-Description: How long something has existed; elapsed time since birth.  
[http://purl.obolibrary.org/obo/NCIT\\_C25150](http://purl.obolibrary.org/obo/NCIT_C25150)

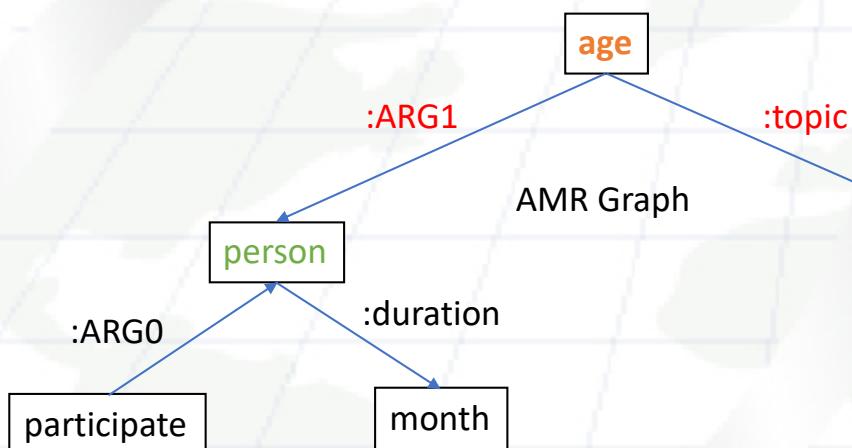
## SDD-Gen Suggestions

Bioportal Search Results



# SDD-Gen

Age in months of the participant at the time of screening.



## AMR entity-related Types

(:ARG0, :ARG1, :ARG2, :ARG3, :conj, :domain, :topic, :location)

[0.6,0.2,...0.9 | 0,...0 | 0.2,0.4,...0.7 | 0,0,...0 | 0,0,...0 | 0,0,...0 | 0,...0 | 0.5,0.7,...1 | 0,0,...0]

GloVe Vectors

Autoencoder

[0.3, 0.6, 0.7, 0.8, 0.2, 0.3, 0.9, 0.6, 0.3, ...0.8]

Lower Dimensional Representation



# Data Integration: sdd2rdf

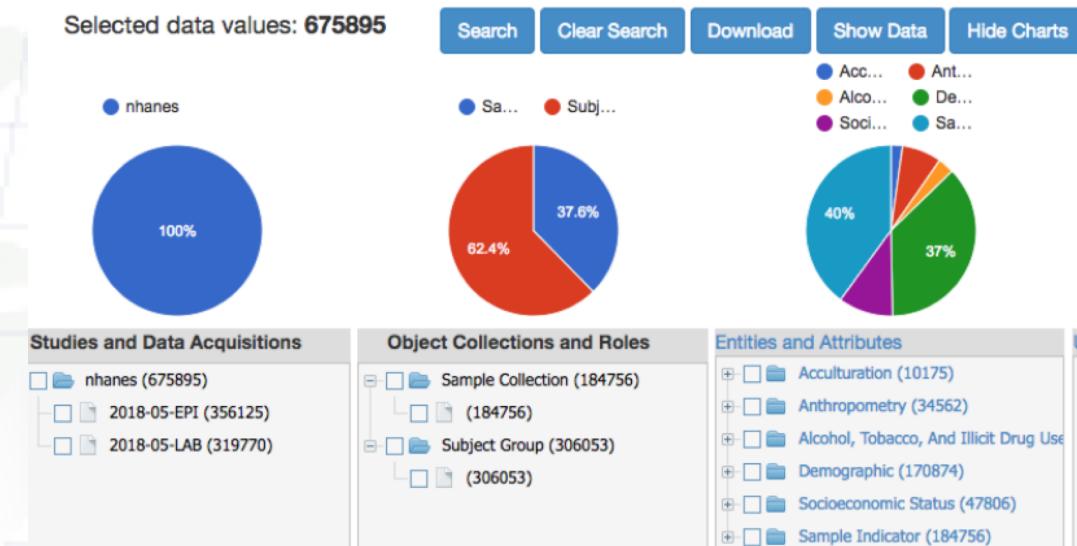
- Lightweight Python library
- Generates knowledge graph fragments as nanopublications
- Applied to a variety of biomedical datasets:
  - NHANES
  - TCGA
  - CIViC
  - SEER
  - MIMIC-III





# Data Integration: HADatAc

- Framework to align data across studies
- Ingests files and message streams
- User Interface:
  - Explore ontologies
  - Ingest datasets
  - Investigate knowledge graphs





# Conclusion

- Semantic Data Dictionaries can be used harmonize diverse datasets
- Our SDD-Editor allows users to semi-automatically generate Semantic Data Dictionaries, thus reducing user burden and broadening the class of users
- HADatAc and sdd2rdf can populate a PHKG using Semantic Data Dictionaries

**Acknowledgements:** This work is partially supported through AFRL 88ABW-2020-0991 and NIH 2U2CES026555-02