



Semantic integration of US Federal nanomaterials data

Holly M. Mortensen

Senior Research Biologist

US EPA/Office of Research and Development

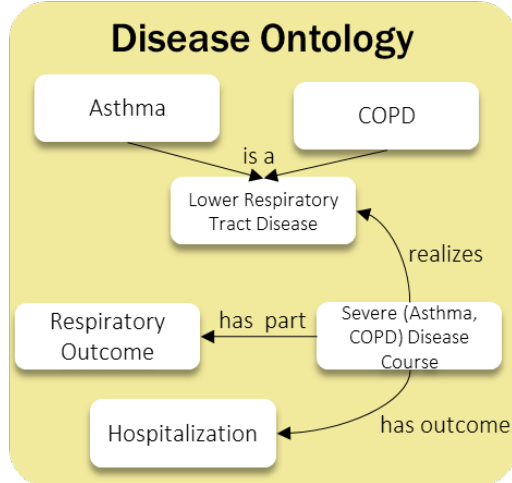
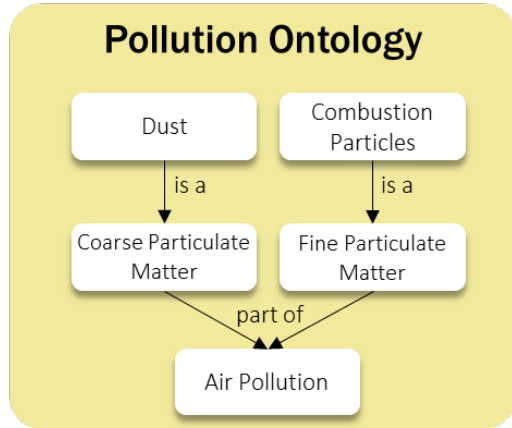
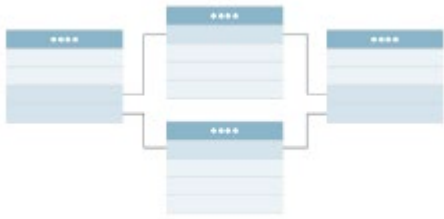
FCSM Research and Policy Conference
Open-Source Software in the Federal Statistical System
Chair: Chris Marcum
University of Maryland, College Park
October 22, 2024, 2pm EST

EPA Disclaimer: The views expressed in this presentation are those of the author(s) and do not necessarily represent the views or policies of the Agency.

Presentation Outline

- Background: ***Knowledge Organization Systems***
 - *Relational vs. Graph representations-examples*
 - **Environmental Health Data Diversity-EHLC**
- EPA Nanomaterials Knowledgebase- ***NaKnowBase***
 - *Motivations-Nomenclature debacle*
 - Proof of Concept-Semantic/Ontology mapping of *NKB* and the EPA *OntoSearcher*
- ***Consortium Effort: NNI NEHI Database Interoperability Group (DIG)***
 - US Federal Agency NanoEHS Consortium Established
- Progress
 - NNI NanoInformatics Conference Nov. 2023; *Conference Proceeding Pub. 2024*
 - 2024 NNI EHS Research Strategy Update (prev. 2011)
 - *Progress and Future directions*

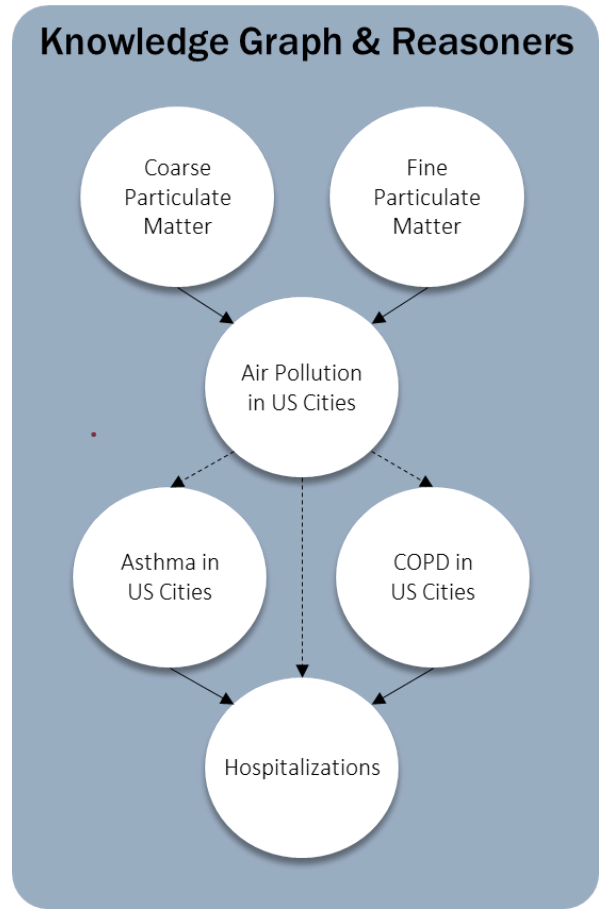
Knowledge Representation



Databases
For Air Quality and Disease Prevalence in US Cities*



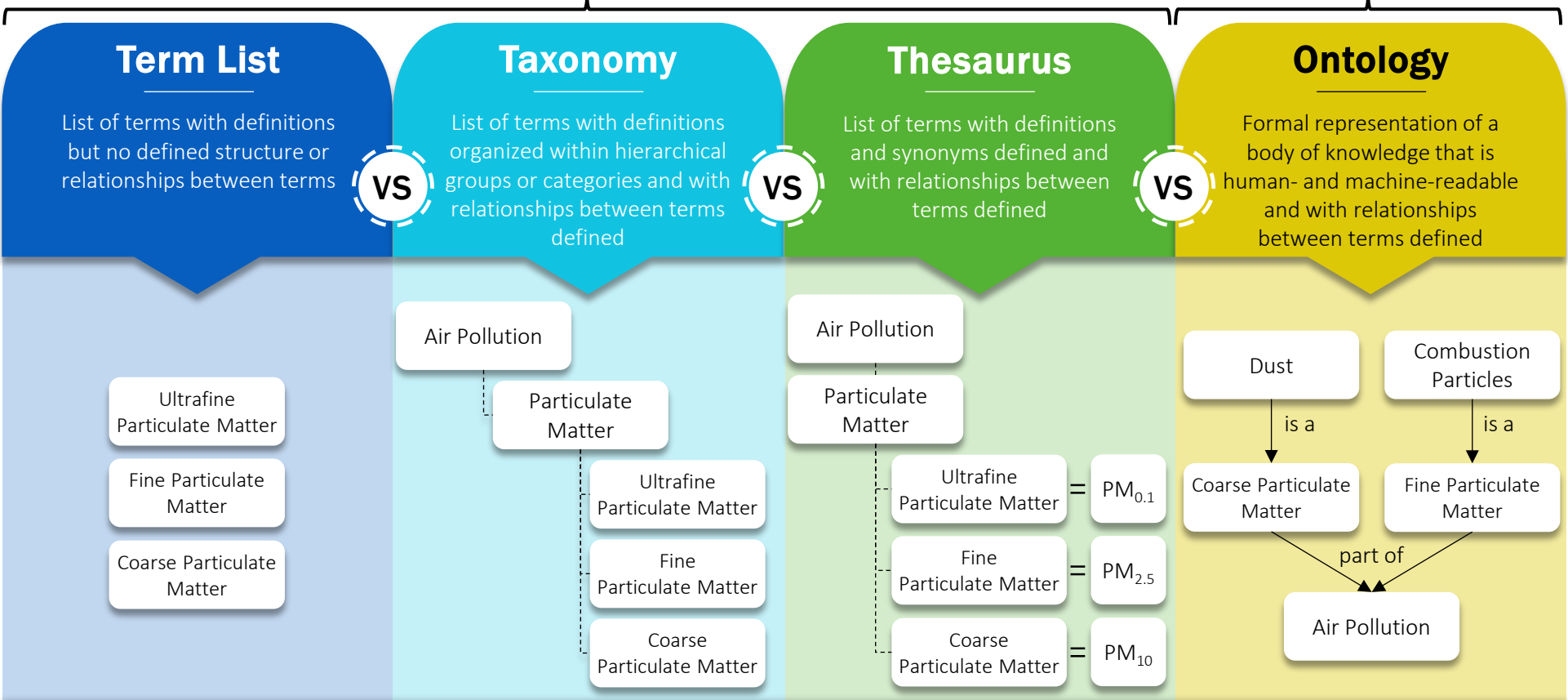
What's driving an increase in hospitalizations?



*Captured using common data elements in a data model with minimal information standards captured using controlled vocabularies

→ Defined Relationship -.-.-> Inferred Relationship

Knowledge Organization
Knowledge Representation



Why Use It?

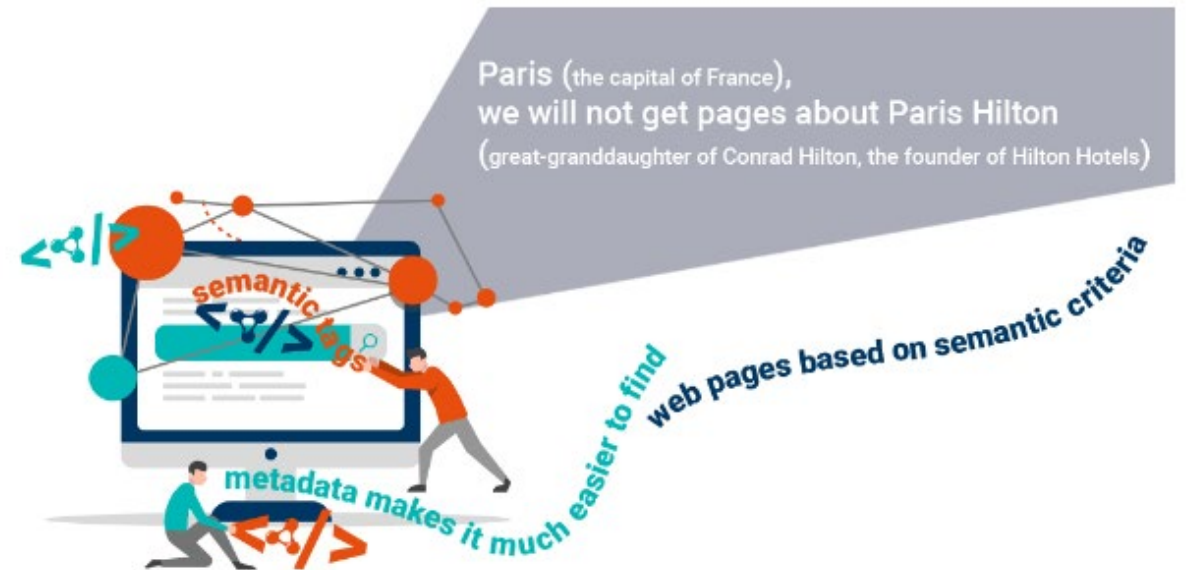
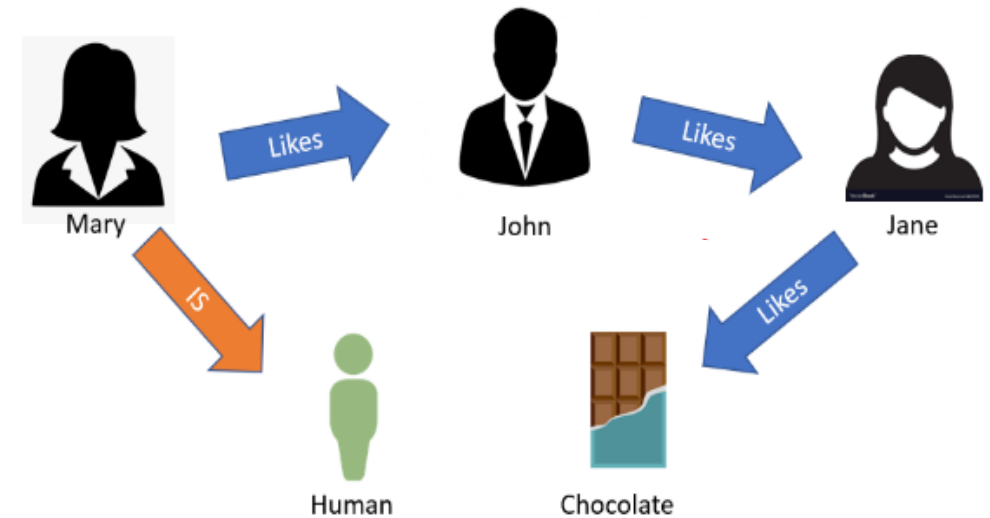
<h4>Consistent</h4> <p>Collaboration and search are easier when everyone is using a set of terms with an agreed upon definition.</p>	<h4>Hierarchical</h4> <p>The hierarchical structure helps us integrate data collected at different levels of granularity and search for data using high level categories.</p>	<h4>Associative</h4> <p>In addition to the hierarchy, the presence of synonyms further improves integration and search over heterogeneous data.</p>	<h4>Inferential</h4> <p>The combination of formal logic and persistent, unique identifiers enables inferencing, makes data computable, and reveals novel connections.</p>
--	---	---	--

Examples

HAWC EHV; LTER Controlled Vocabulary	NCBITaxon	NCIThesaurus; MeSH	OBO Foundry Ontologies; UBERON
--------------------------------------	-----------	--------------------	--------------------------------

Some Terminology...

- Semantic mapping is a way of representing information (concepts or data) as a **graph**
- Resource Description Framework (RDF) is a **directed graph** AND a **data model** for exchanging information on the web
 - triples= subject, predicate, and object
- Is RDF appropriate for nanoEHS data?
 - "**Nomenclature debacle**" - lack of consistent nomenclature across sources exacerbates integration
- **Concept**-By adding in the metadata component, semantic technologies can address data heterogeneity and interpolation issues



Environmental Health Data Diversity: *subfield contributors, language reporting standards, and actors and stakeholders*



The Environmental Health Language Collaborative

Harmonizing Data. Connecting Knowledge. Improving Health.



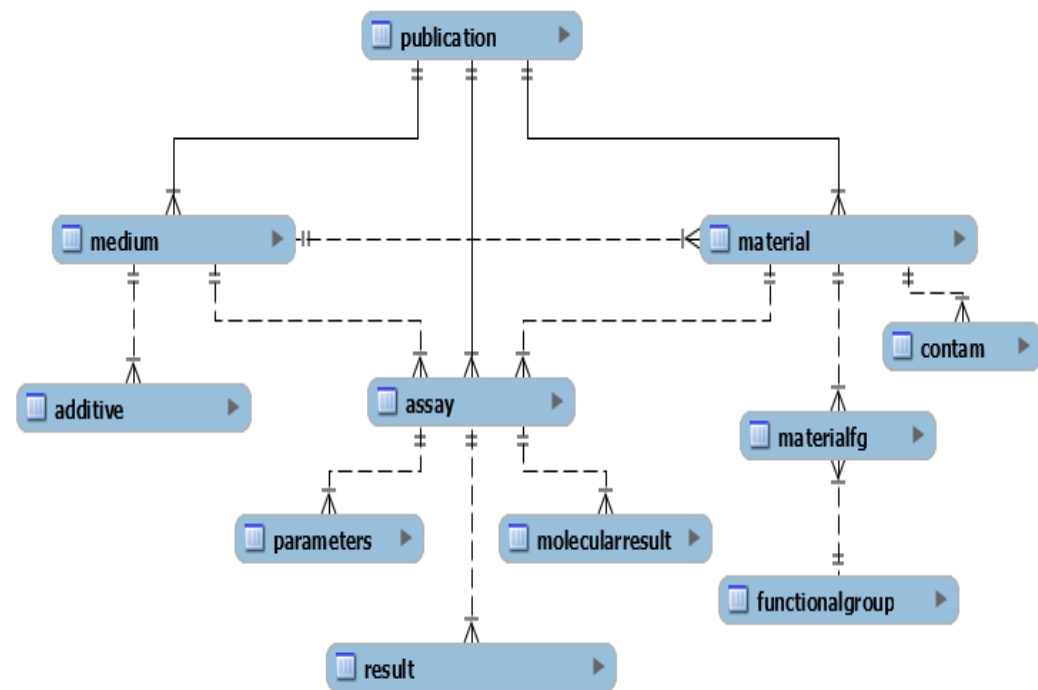


EPA NaKnowBase

A Curated, Relational Database detailing physical-chemical properties of EPA emerging materials research

EPA OntoSearcher

Automated assignment of ontology terms and graph creation



Relational SQL database

Curated Data Fields:

- Publication (DOI, author etc.) **>120 peer-reviewed manuscripts (2012-2019)**
- Materials (**>70 unique NM**)
 - Physical/Chemical properties
 - Capping materials
 - Media
 - Contaminants
- Assays (**>160 named assays**)
 - Parameters Measured
- Results

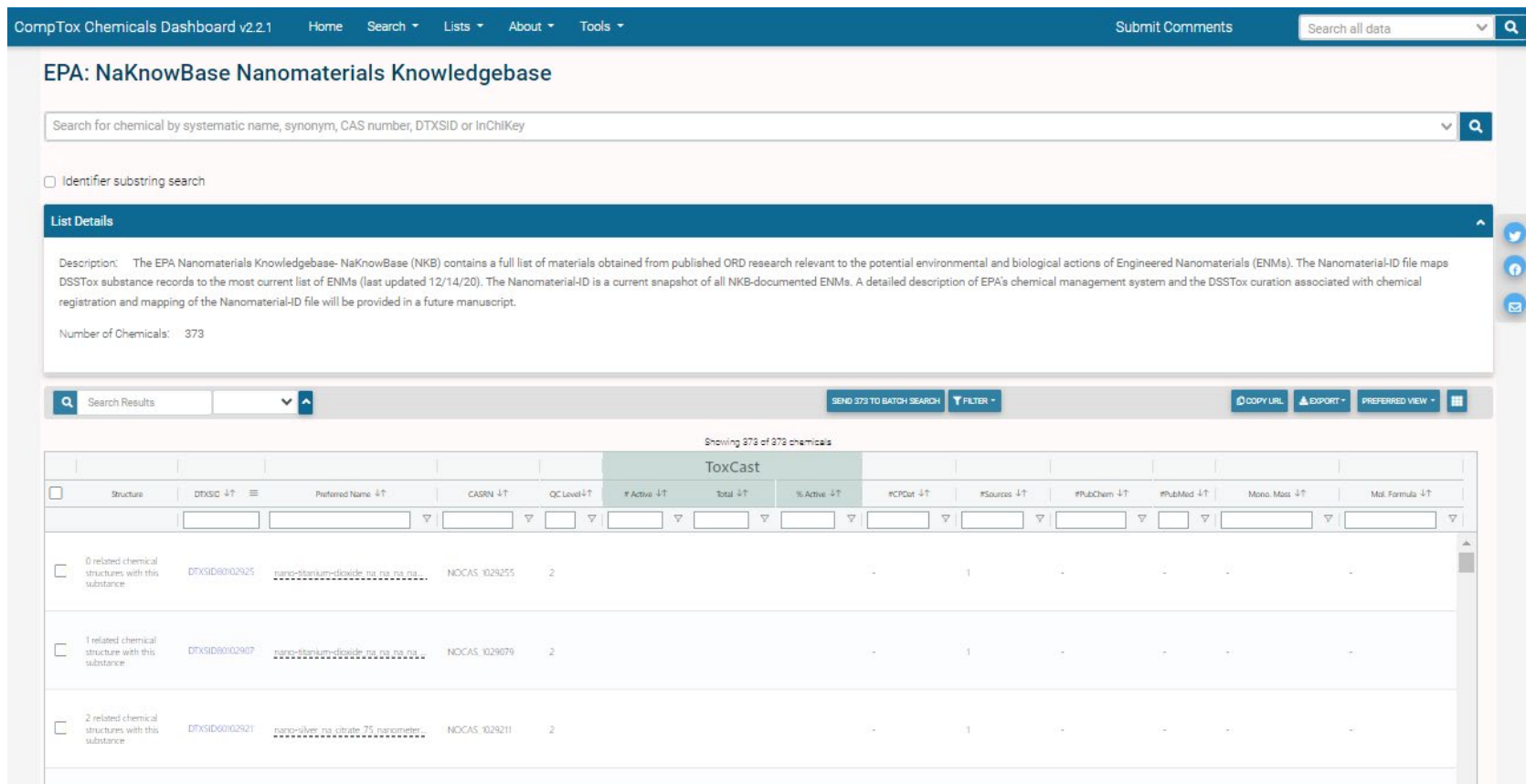
Boyes, W.K., Beach, B., Gayle Chan, G., Thornton, B.M., Harten, P., Mortensen, H.M. (2022) An EPA database on the effects of engineered nanomaterials-NaKnowBase. *Nature Sci Data* 9, 12. <https://doi.org/10.1038/s41597-021-01098-0>.

Natural Language Processing (NLP) for Computable and Interoperable Descriptions of EPA nanomaterials

nano-SiO₂ NA NA NA nm ENPRA A

This nano scaled material is composed of a SiO₂ core with no information on the surface coating or capping agent. There is no data on manufacturer reported particle size (diameter). This material was obtained from ENPRA as sample A of the material of the same core, coating/capping, diameter and source information.

NKB nanomaterials on the EPA CompTox Chemistry Dashboard



The screenshot shows the EPA CompTox Chemistry Dashboard interface. At the top, there is a navigation bar with 'Home', 'Search', 'Lists', 'About', and 'Tools' menus. A search bar is present with the text 'Search all data'. Below this, the page title is 'EPA: NaKnowBase Nanomaterials Knowledgebase'. A search input field contains the text 'Search for chemical by systematic name, synonym, CAS number, DTXSID or InChIKey'. A description of the knowledgebase is provided, stating that it contains a full list of materials from published ORD research on Engineered Nanomaterials (ENMs). The number of chemicals is listed as 373. A table of results is shown, with columns for Structure, DTXSID, Preferred Name, CASRN, QC Level, # Active, Total, % Active, #CPDat, #Sources, #PubChem, #PubMed, Mono. Mass, and Mol. Formula. The table displays three rows of chemical entries.

Structure	DTXSID	Preferred Name	CASRN	QC Level	# Active	Total	% Active	#CPDat	#Sources	#PubChem	#PubMed	Mono. Mass	Mol. Formula
0 related chemical structures with this substance	DTXSID00102925	nano-barium-dioxide	NOCAS: 102925	2	1	1	100%	1	1	1	1		
1 related chemical structure with this substance	DTXSID00102907	nano-barium-dioxide	NOCAS: 102907	2	1	1	100%	1	1	1	1		
2 related chemical structures with this substance	DTXSID00102921	nano-silver-nitrate 75 nanometer	NOCAS: 102921	2	1	1	100%	1	1	1	1		

The Nanomaterial-ID file maps DSSTox substance records

- **373 ENMs mapped to DSSTox IDs**

EPA's **DSSTox (Distributed Structure-Searchable Toxicity)** database contains curated chemical substances mapped to chemical identifiers (i.e., chemical synonyms, systematic names, CAS Registry Numbers and others) and, where appropriate, chemical structure representations.

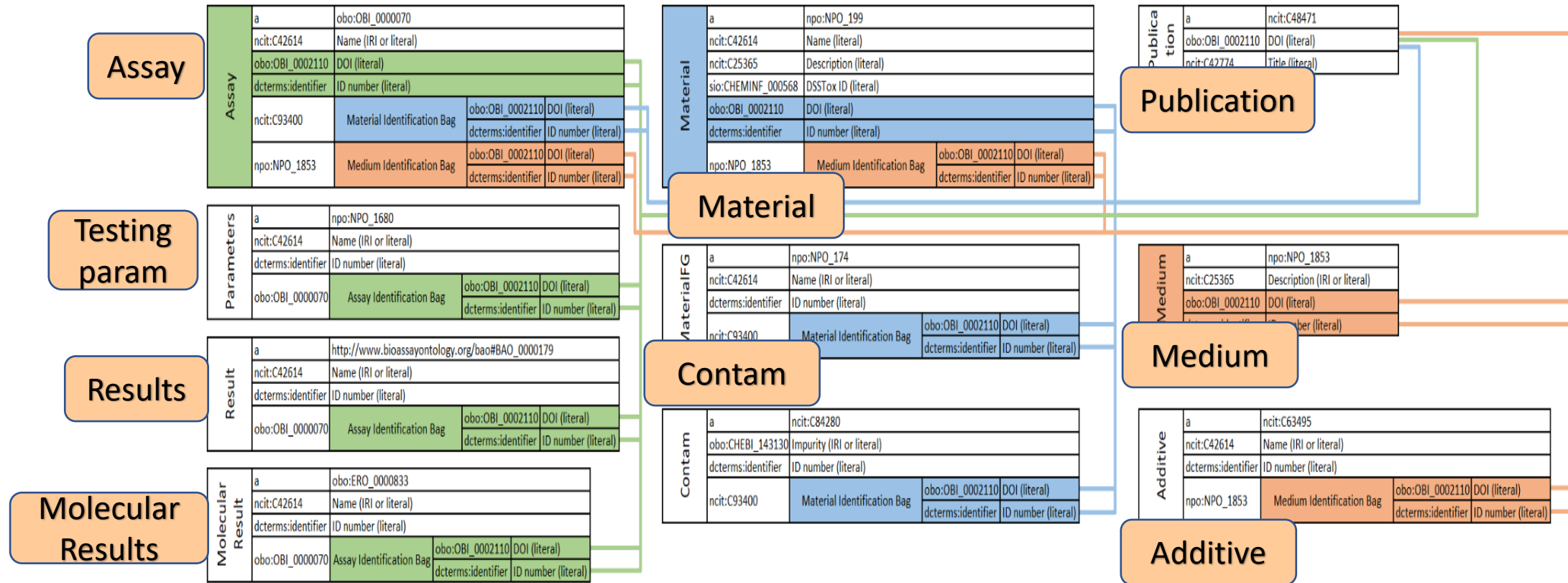
Naknowbase About How to Cite Tools ▾

1 Material ⇅ 2 CoreComposition ⇅ 3 = ⇅ 4 value ⇅ X

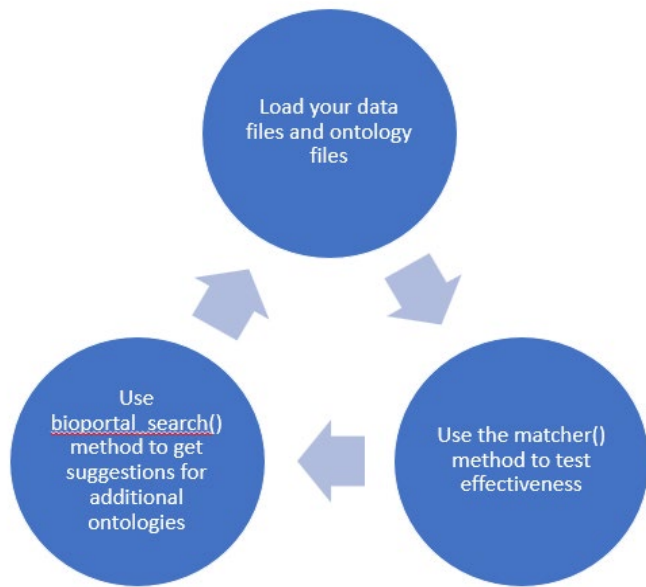
1. List of tables in NKB
2. List of fields in the selected table.
3. Comparison symbols for numeric fields.
4. List of values for non-numeric fields, or a text box for numeric fields.

+ - Submit Restart

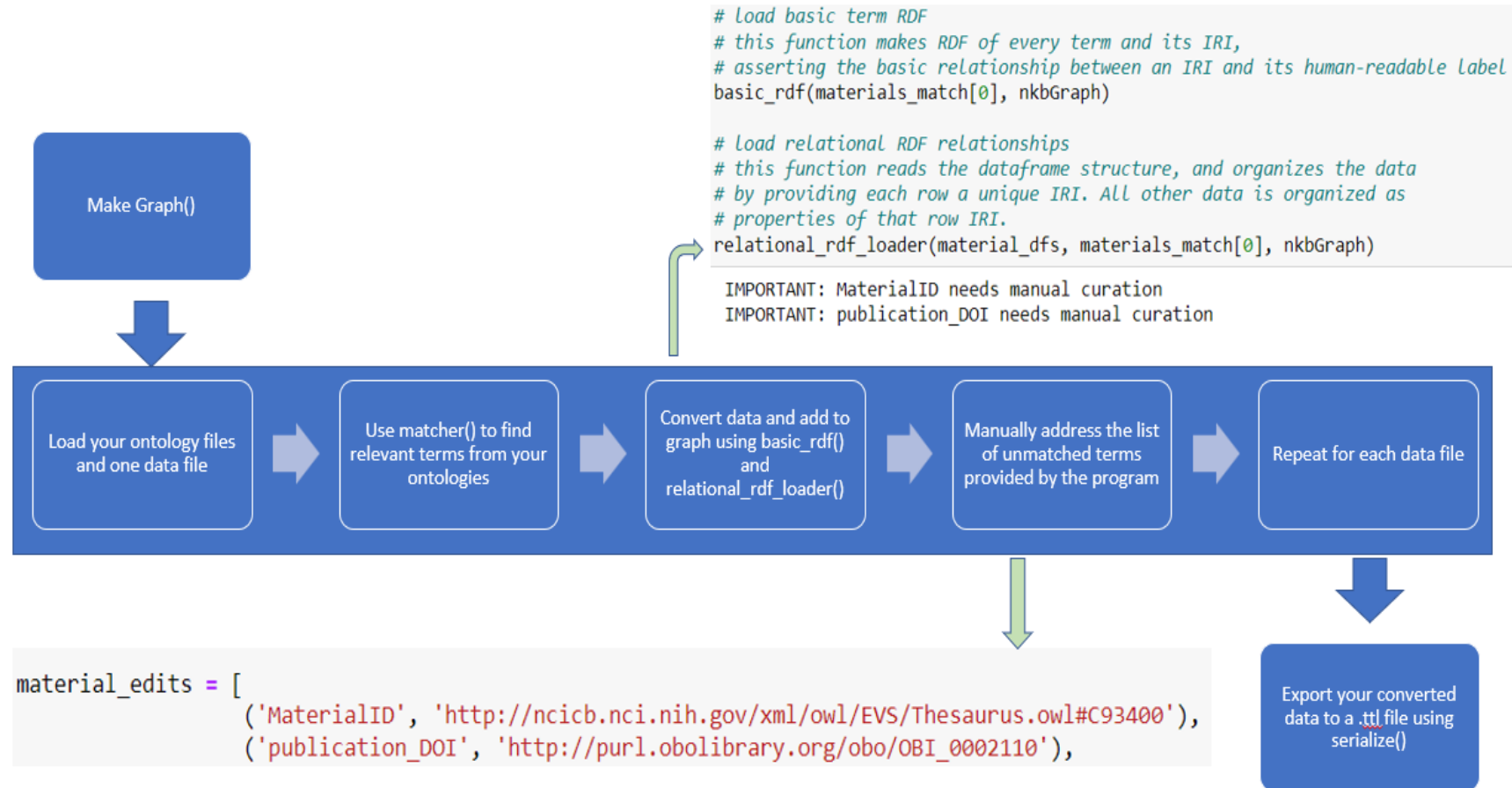
Submit a Query... Internet Subdomain Request Initiated---HOLD:
naknowbase.epa.gov



EPA OntoSearcher for Automated Term Mapping and Software-guided graph creation



- Automated mapping with match score to Assist in manual curation
- Convert data into RDF
- .csv input to .ttl output



Federated SPARQL Query – AOP-DB pathways by AOP-DB gene and NKB material

```
#determine which NKB material/AOPDB gene combinations are associated with the most pathways
aopdb_fed = """
SELECT distinct ?geneID (COUNT(?pathwayname) as ?p) ?DOI ?material ?DTXSID
WHERE {
  SERVICE <https://aopdb.rdf.bigcat-bioinformatics.org/sparql> {
    ?chemgeneassoc a <http://semanticscience.org/resource/SIO_001257>.
    ?chemgeneassoc <http://semanticscience.org/resource/CHEMINF_000446> ?CAS.
    BIND(URI(CONCAT("http://identifiers.org/cas:", STRAFTER(str(?CAS), "https://identifiers.org/cas:"))) as ?CAS2).
    ?chemgeneassoc <http://edamontology.org/data_1027> ?geneID.
    ?pathway <http://edamontology.org/data_1027> ?geneID.
    ?pathway a <http://purl.obolibrary.org/obo/PW_0000001>.
    ?pathway dc:title ?pathwayname.
  }
  ?material sio:CHEMINF_000446 ?CAS2.
  ?material a npo:NPO_199.
  ?material npo:NPO_1808 ?core.
  ?material obo:OBI_0002110 ?DOI.
  ?material sio:CHEMINF_000568 ?DTXSID.
}
GROUP BY ?geneID
ORDER BY DESC(?p)
LIMIT 5
"""
qres = g.query(aopdb_fed)
for row in qres:
    print(f"{row.DTxsID} of {row.DOI} | gene {row.geneID} | \n# pathways: {row.p}")

http://identifiers.org/comptox/DTXSID501028989 of 10.1186/s12951-014-0047-3 | gene https://identifiers.org/ncbigene/196 |
# pathways: 560
http://identifiers.org/comptox/DTXSID501028989 of 10.1186/s12951-014-0047-3 | gene https://identifiers.org/ncbigene/208 |
# pathways: 283
http://identifiers.org/comptox/DTXSID501028989 of 10.1186/s12951-014-0047-3 | gene https://identifiers.org/ncbigene/154 |
# pathways: 196
http://identifiers.org/comptox/DTXSID301028969 of 10.1016/j.watres.2012.12.041 | gene https://identifiers.org/ncbigene/180359 |
# pathways: 56
http://identifiers.org/comptox/DTXSID501028989 of 10.1186/s12951-014-0047-3 | gene https://identifiers.org/ncbigene/948 |
# pathways: 53
```

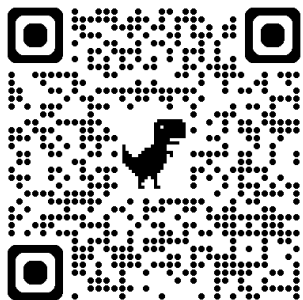
Federated Query calling for nanomaterials and corresponding gene targets that hit the greatest number of pathways



EPA Interoperability tools and Relational DB products



DATA.GOV: NKB Data Catalog
<https://catalog.data.gov/dataset/naknowbase-interoperability-tools>



The screenshot shows the DATA.GOV website interface. At the top, there are navigation links for DATA, REPORTS, OPEN GOVERNMENT, and CONTACT. Below this is a blue header with 'DATA CATALOG' and buttons for 'Datasets' and 'Organizations'. The main content area displays the dataset 'NaKnowBase Interoperability Tools' from the U.S. Environmental Protection Agency. It includes a description of the dataset, metadata (updated September 17, 2023), and a detailed text description of the dataset's contents and access methods. There are also sections for 'Access & Use Information' (Public, License) and 'Downloads & Resources' with links to the dataset page and training materials.



US Federal Agency NanoEHS Consortium

**Informatics plan described in
2030 US EU Roadmap and 2024 NNI EHS Research Strategy**

Introductory Results and Project plan

US Federal Agency NanoEHS Consortium Established

Welcome Home Find Help Contact Us Search the MAX Community All Welcome, Holly

MAX.gov Shared Services will be sunseting in December 2023. [Click here](#) for an important update from November 2023.

PERMISSIONS OPEN-EXECUTIVE BRANCH Edit Add Favorites Watchers Share Actions

Pages / ... / nanoEHS Databases and Related Resources

NANOEHS DATA REPOSITORY

Created by Rhema Bjorkland (ARC), last modified on Jun 25, 2021

Please contact Rhema Bjorkland (rbjorkland@nnco.nano.gov) to request permission to upload you data.

June 16, 2021: We encourage you to add a readme or metadata file when you upload your data or information to the page. Please include your name, affiliation, and contact information as well.

June 25, 2021: Adding multiple documents under a child page is a helpful organizational approach. Use the green "Add child page" link to create a page, giving the page your agency name. Add documents under the child page using the "Attachments" link.

No labels

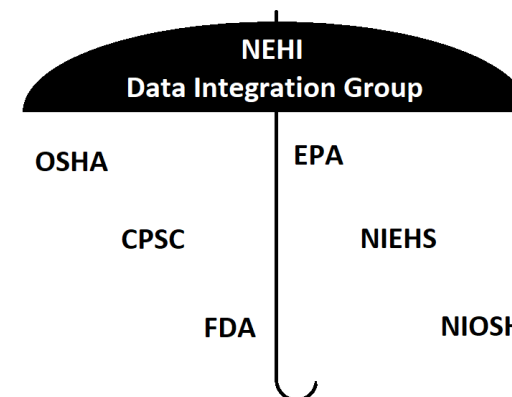
Child Pages (1) Add Child from Template Add Child Page Attachments (2) Sort Show Details Advanced Add Attachment(s)

NIOSH Data Templates

OSHA_Nanomaterial Database_sample_06282021.xlsx (33 KB, v.1)
Last edited by: Rhema Bjorkland (ARC) on Jun 28, 2021 at 09:51 AM
 OSHA Sample. Janet Carter is the POC.
No labels

nanowbase_schema_10-22-2020.png (141 KB, v.1)
Last edited by: Holly Mortensen (EPA) on Jun 26, 2021 at 11:07 AM
 NKB Schema
No labels

Comments (0) Add Comment



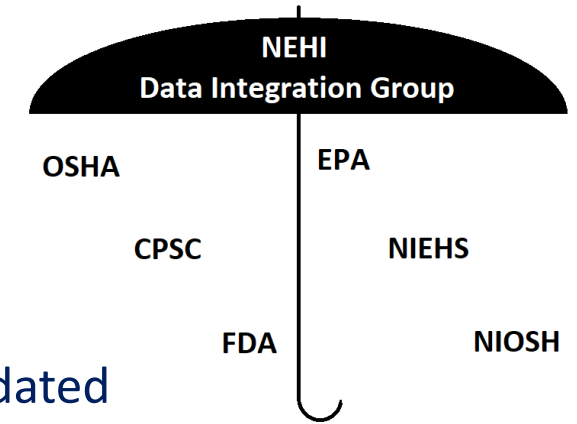
- **2022 NEHI-DIG Consortium Formed**
- Initial Data repository established (MAX.gov)
- Participating Federal Agencies: EPA, NIOSH, OSHA, CPSC, FDA, NIH

<https://www.nano.gov/NNINanoinformaticsConference>

**Described in the upcoming NNI EHS Strategy update 2023: Informatics and Modeling. Federal Register/Vol. 89, No. 115/Thursday, June 13, 2024/*

Motivations for forming the DIG Consortium

- Shared vision of Federal partners – team science, and data sharing
 - *breaking down data silos!*
- **Mechanistic Interaction** of engineered nanomaterials (ENMs) is not yet fully elucidated
- Lack of information on **toxic relevance** (e.g. which disease outcomes are relevant for ENMs?)
- *Computationally structured, semantic annotation* can improve our ability to understand this biology
 - Promoting FAIR (Findable, Accessible, Interoperable, and Reuseable) data management and sharing principles for ENM would simplify data integration with other knowledge systems
 - EPA NKB as proof of concept- *We can do this!*
 - Stay current with EU progress in this area (??)



NIOSH: EPA processing with OntoSearcher

EPA OntoSearcher: CSV to RDF Conversion

NIOSH Dataset

This document uses EPA's OntoSearcher application to convert multiple CSVs, derived from an Excel workbook of nanomaterial research data provided by NIOSH, into **Resource Description Framework (RDF)**. RDF is a data format which uses unique web addresses, called **Internationalized Resource Identifiers (IRIs)**, to identify pieces of unique information. Associating data with these unique identifiers and publishing that data in RDF format allows for any data regarding the same entity (that shares an IRI) to be interoperable.

EPA OntoSearcher is a prototype application developed at the **Dr. Holly Mortensen** lab at **EPA ORD CPHEA** to expedite the conversion of relational data into RDF. The application provides functions for importing CSV data, importing ontology and RDF data, search algorithm functions to compile a dictionary of IRI's for csv terms, and functions that build RDF from csv data and term-IRI associations.

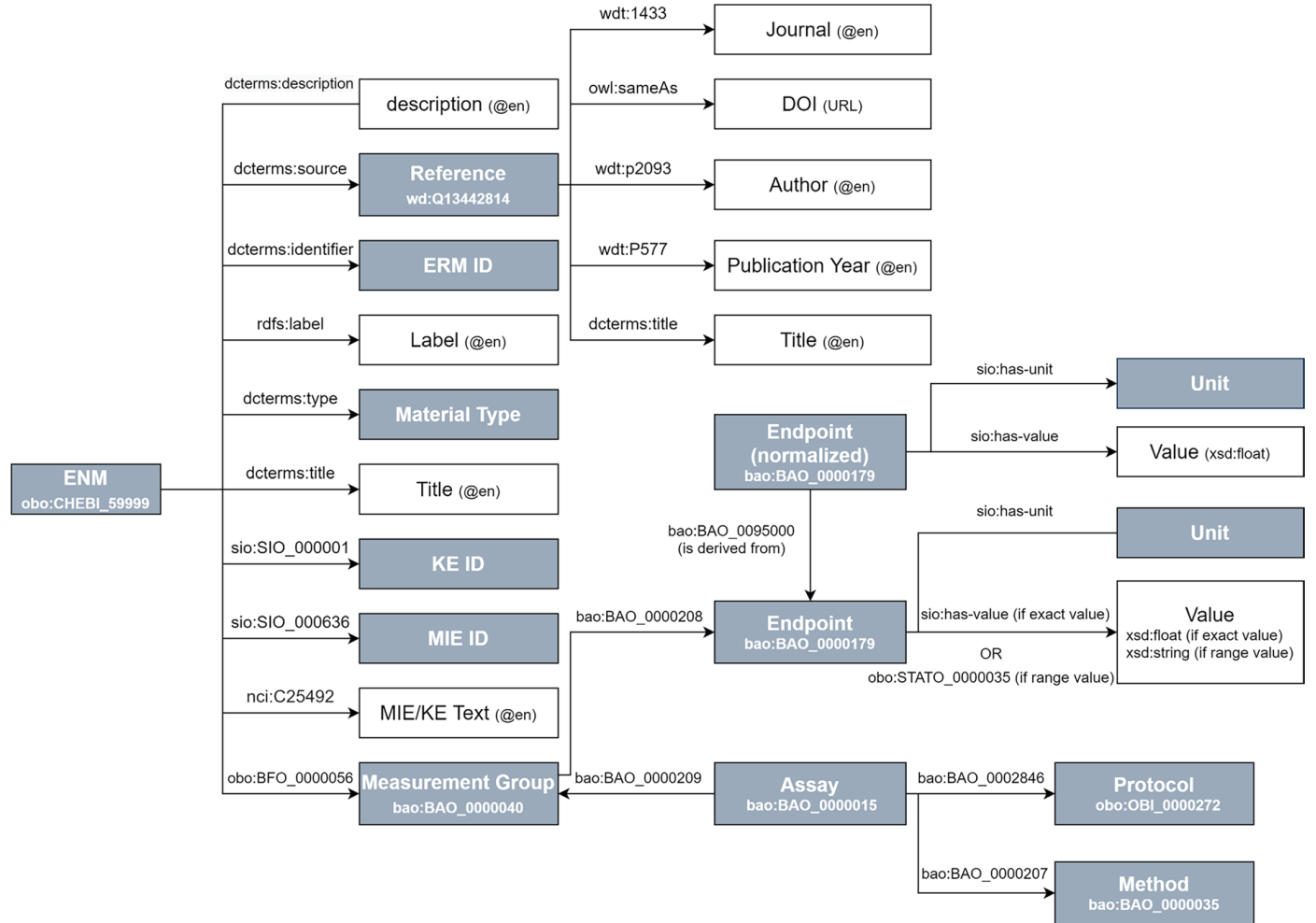
This document will showcase all of this functionality, as well as how to query RDF data using SPARQL, the RDF query language.

Why are we here?

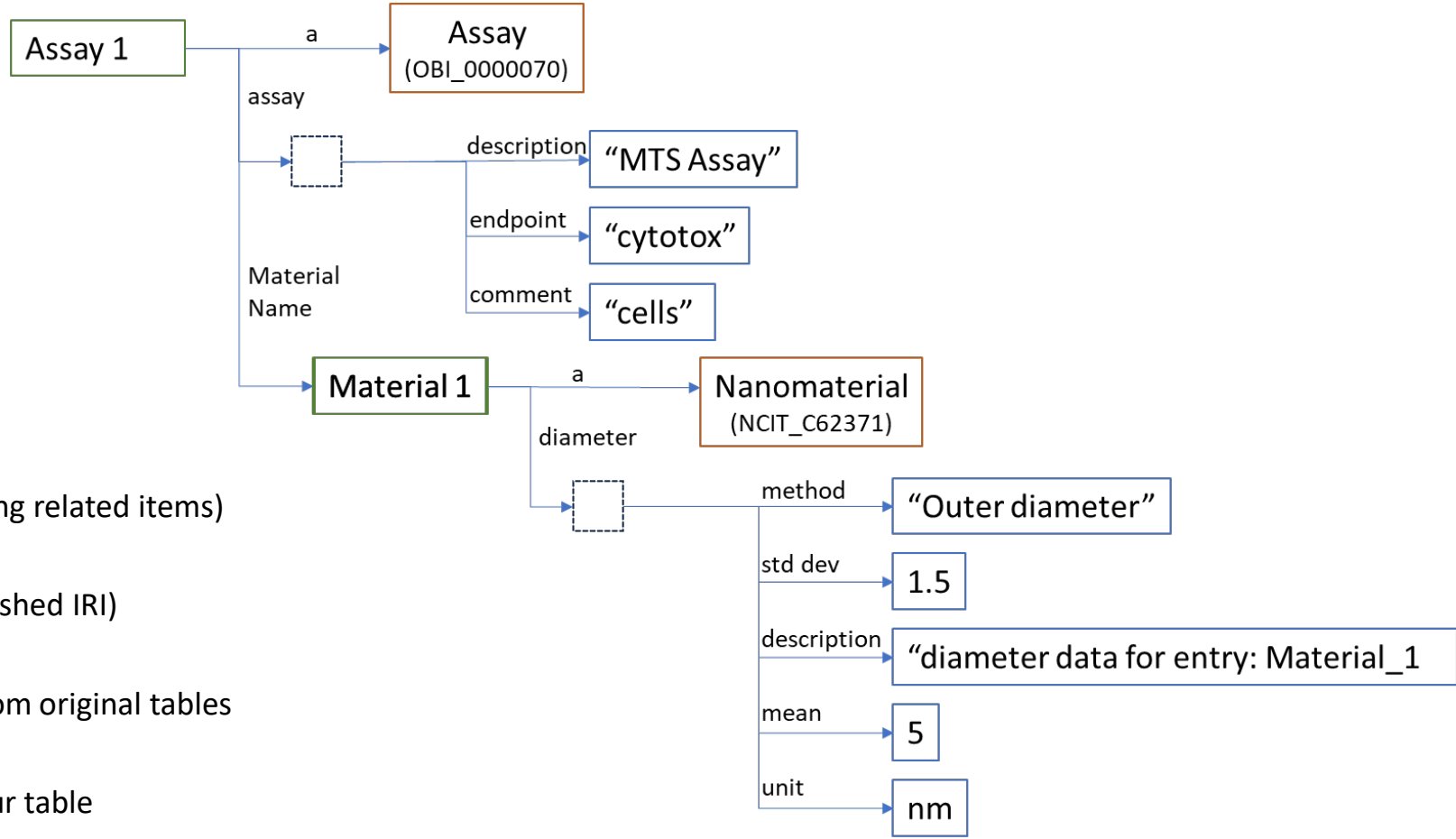
to answer *three questions*

- How should we **format our relational data** to make interoperability easier?
- What is **RDF/OWL** and **what utility does it have** for my needs?
- **How can I convert my data** into RDF/OWL (without breaking a sweat)?


```
# import EPA OntoSearcher modules, and other packages
from onto import ontolister, ontocontext
from csv_importer import load_data
from find import matcher
from onto_api import bioportal_search, unpack_superclass
from onto_api import bioportal_sample, dict_samp, bio_summary
from rdf_print import table_from_file, term_editor, term_lookup
from rdf_print import basic_rdf, relational_rdf_loader
from rdf_print import primenode, node_one, node_two, multi_editor
```





NIOSH: EPA manual interrogation





Key

 Blank Node (useful for grouping related items)

 Predicate (~column w/ established IRI)

 Starter node/row indicator from original tables

 Literal value, the data from our table

 Object with established IRI (extra metadata)

NNI Nanoinformatics Conference



November 15, 2023
In-person, Washington, D.C.

Computational Toxicology 30 (2024) 100316

Contents lists available at ScienceDirect

Computational Toxicology

journal homepage: www.sciencedirect.com/journal/computational-toxicology



Short Communication

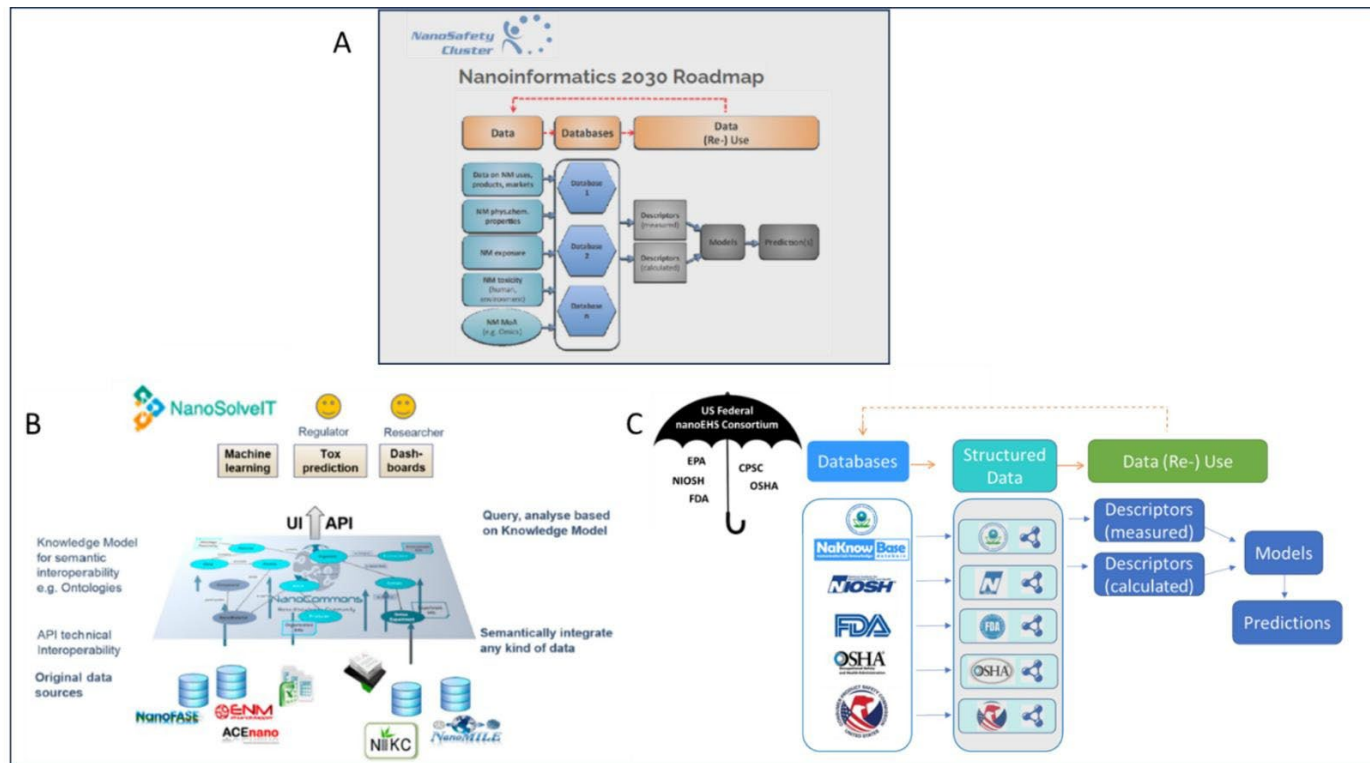
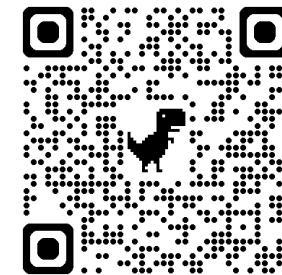
NNI nanoinformatics conference 2023: Movement toward a common infrastructure for federal nanoEHS data computational toxicology: Short communication

Holly M. Mortensen^{a,*}, Jaleesia D. Amos^{b,2}, Thomas E. Exner^{c,3}, Kenneth Flores^{d,4}, Stacey Harper^{e,5}, Annie M. Jarabek^{f,6}, Fred Klaessig^{g,7}, Vladimir Lobaskin^{h,8}, Iseult Lynch^{i,9}, Christopher S. Marcum^{j,10}, Marvin Martens^{k,11}, Branden Brough^l, Quinn Spadola^{m,12}, Rhema Bjorkland^{m,13}

^a Public Health and Integrated Toxicology Division, Center for Public Health and Environmental Assessment, Office of Research and Development, US Environmental Protection Agency, Research Triangle Park, NC, United States
^b Civil & Environmental Engineering, Duke University, Durham, NC 27708, United States
^c Seven Post Nine GmbH, Reibacker 68, 79650 Schopfheim, Germany
^d NSF MPS Ascend Fellow, Arizona State University, Tempe AZ 85287, United States
^e Department of Environmental and Molecular Toxicology and the School of Chemical, Biological and Environmental Engineering, Oregon State University, Corvallis, OR 97331, United States
^f Health and Environment Effects Assessment Division, Center for Public Health and Environmental Assessment, Office of Research and Development, US Environmental Protection Agency, Research Triangle Park, NC 27711, United States
^g Pennsylvania Bio Nano Systems, PA, United States
^h School of Physics, University College Dublin, Dublin 4, Ireland
ⁱ University of Birmingham, Birmingham, UK
^j Office of the Chief Statistician of the United States, Office of Management and Budget, Washington, DC 20503, United States
^k Department of Bioinformatics - BIGaT, NUTRIM, Maastricht University, Maastricht, the Netherlands
^l National Nanotechnology Coordination Office, Alexandria, VA, United States
^m Contractor to the National Nanotechnology Coordination Office, Alexandria, VA, United States

NNI nanoinformatics conference 2023

<https://doi.org/10.1016/j.comtox.2024.100316>



EU US Roadmap Nanoinformatics 2030

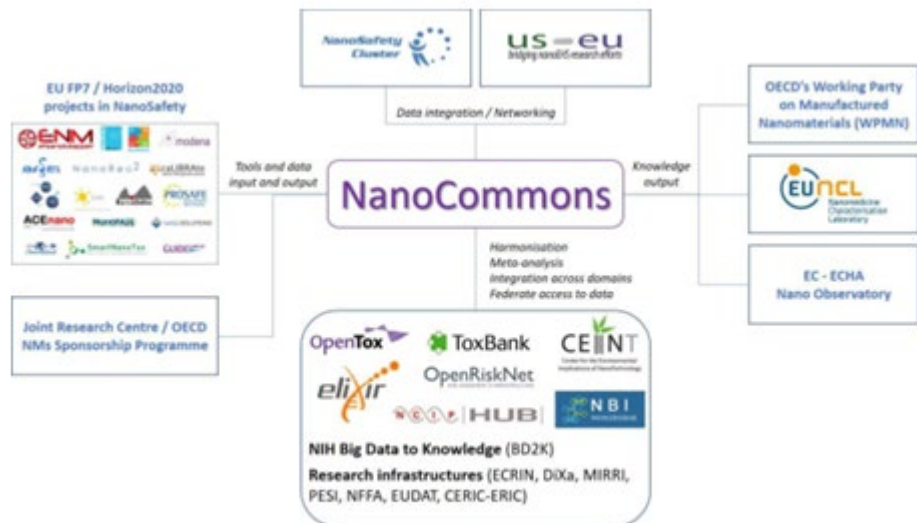
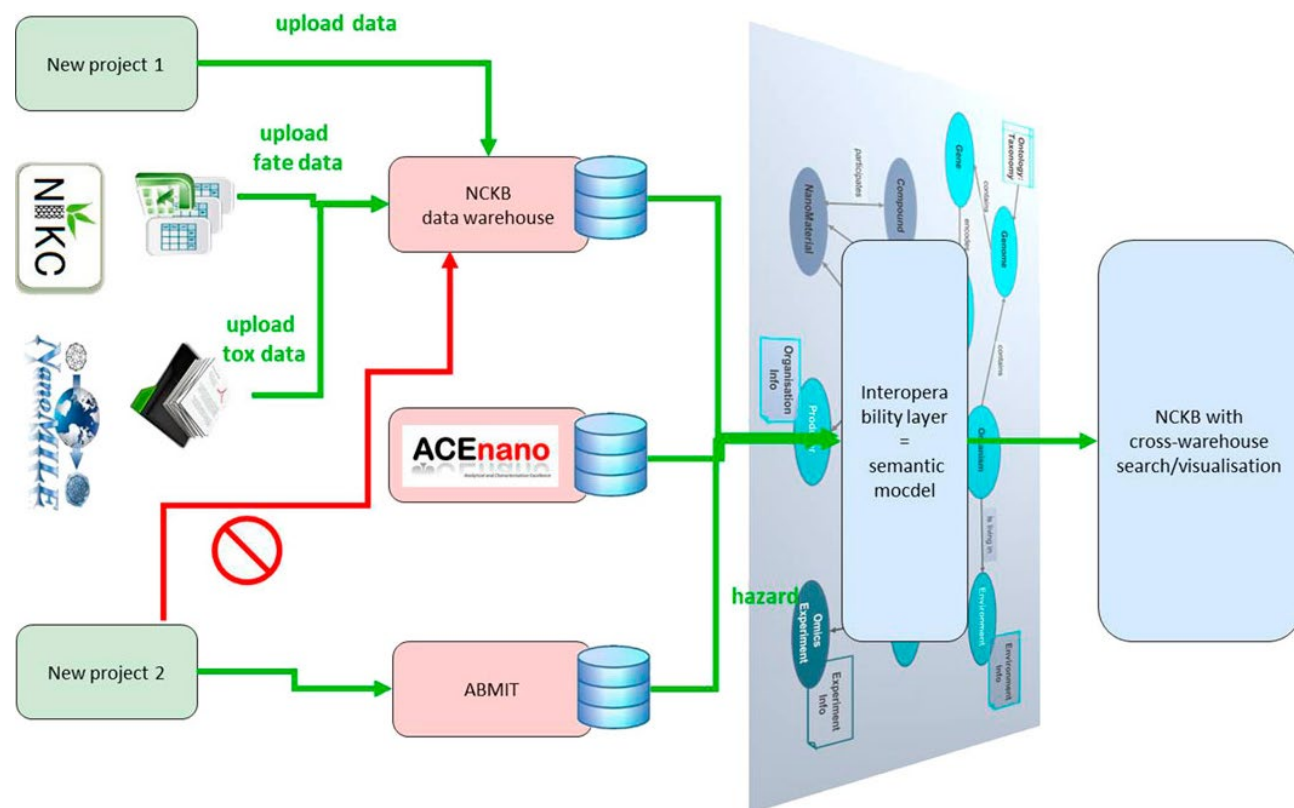


Figure 16: Schematic illustration of the positioning of NanoCommons and how it will provide an integrating platform for the nanosafety knowledge community in Europe and internationally.



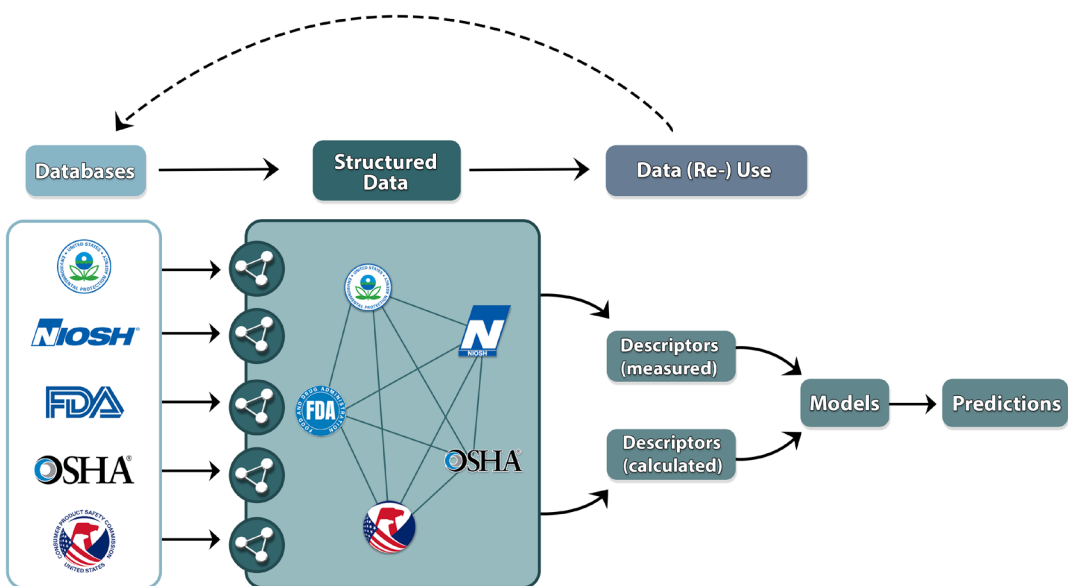
Maier, et al. Front. Phys., 12 November 2023
<https://doi.org/10.3389/fphy.2023.1271842>



2024 NNI EHS Research Strategy

NATIONAL NANOTECHNOLOGY INITIATIVE
ENVIRONMENTAL, HEALTH, AND SAFETY
RESEARCH STRATEGY: 2024 UPDATE

DRAFT FOR PUBLIC COMMENT, 06/10/24



- Identifies Nanoinformatics as a cross-cutting theme
- Needs
 - Expand and Strengthen the Collaborative Informatics Infrastructure
 - Boost informatics and data infrastructure for robust risk assessment and decision-making
 - Ensure alignment with FAIR and TRUST principles:

Requires the development of standardized protocols for data reporting and sharing, ensuring that data generated from nanotechnology research is easily findable and accessible to a broad range of stakeholders

2024 Project proposed



ATTACHMENT I- Project Topic

Federal Government Data Integration and Usage Platform for Emerging and Nanomaterial Environmental and Health Safety

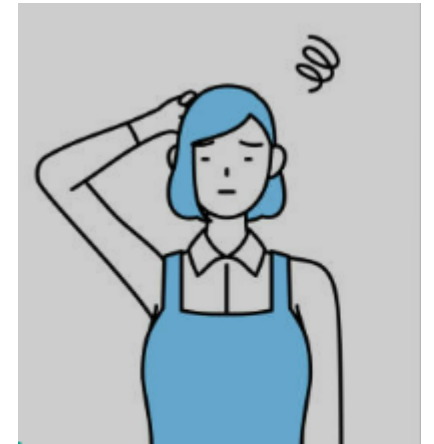
1. A signed **MOU or other agreement vehicle** establishing data sharing, collection, usage.
2. A **web tool** to support partner data sharing and accessibility, as well as semantic interoperability using the *proof-of-concept* tools.
3. A **data model** for contributing federal partner nanoEHS data.
4. A prototype data usage dashboard or other presentation and tools.
5. **Communication Plan for** implementation and deployment across federal agencies.
6. Training data (if applicable), and any other data created under this award.

Future work and Questions/Concerns

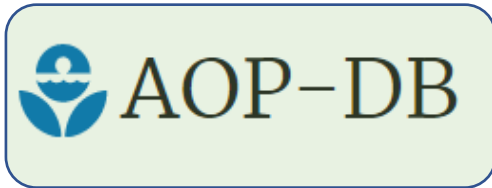
- **Staffing and project support**-NNI interns assigned to expedite processing (Prakash Pranav -UC San Diego)
 - (Post-OntoSearcher) Manual Curation – **Agency Data owners are the lag**
- Are one-to-one reuse of existing ontologies the way to go?
 - Different ontologies use different structures- **Combining terms reused from other ontologies with new domain-specific terms proves consistently problematic**
 - Large number of terms map into the ontology! But...the new ontology has a **specific application domain**
 - Precisely the domain-specific terms that are most relevant, but not available somewhere else.
- **What is needed: Agency coordination and needed expertise**

Planned virtual all-day event

EHLC Use case focusing on AOP-biomedical entity mapping (Early 2025) –contact mortensen.holly@epa.gov



LINKS to EPA projects



- EPA NaKnowBase and related tools
<https://catalog.data.gov/dataset/naknowbase-interoperability-tools>
- AOP-DB web user interface <https://aopdb.epa.gov/>
(permanently decommissioned as of 2024)
- The AOP-DB v.2 stressor linkages are provided through the CompTox Chemicals Dashboard-Number of Chemicals: 349:
https://comptox.epa.gov/dashboard/chemical_lists/AOPSTRESSORS
- AOP-DB SPARQL Endpoint:
<https://github.com/BiGCAT-UM/AOP-DB-RDF>
https://aopwiki.org/info_pages/8

Acknowledgements

EPA

Will Boyes

Antony Williams

EPA-NPD

Annette Guiseppi-Elie

Kathy Dionisio

EPA Students/ Contractors

Bradley Beach

Weston Slaughter

Jonathan Senn

Bradley Sutliff

US Federal Partners

Treye Thomas **CPSC**

Joanna Matheson **CPSC**

Janet Carter **OSHA**

Jay Vietas **CDC/NIOSH**

Kuempel, Eileen **CDC/NIOSH**

Nathan Drew **CDC/NIOSH**

Anil Patri **FDA**

Duke-CIENT/InFRAMES

Mark Weisner

Jaleesia Amos

International Partners

Thomas Exner

Egon Willighagen

Marvin Martens

Andrea Haase

Penny Nymark

NNI-NEHI

Branden Brough

Quinn Spadola

Rhema Bjorkland

Geoff Holdridge

EPA National Program Support: Chemical Safety and Sustainability, RA3: Emerging Materials and Technologies, CSS 403.1: Evaluate environmental impacts of emerging materials on humans and ecological species, Product 403.1.4: Improved NaKnowBase Data Integration to Meet Program and Partner Needs

AND RA8: Informatics, Synthesis, and Integration: CSS.408.2 - Knowledge delivery and interoperability in support of chemical safety decisions, Product 408.2.22 Development of infrastructure support for EMT: EPA NaKnowBase