# *What End Users Need from the Ontology Community*

- *Experience from NCPI, FDA, and COVID-19 Ontologies Harmonization effort.*

*Asiyah Yu Lin, DATA scholar @ NHGRI/NIH*

WSBO2021

*July 15, 2021*

# End users: which ontology do I use?

- End user :
  - developer, domain expert, project manager
  - non-ontologist, not involved in ontology development
- <u>MONDO, DO, HPO? Or SNOMED CT, MeSH</u>…
- PubChem or ChEBI?
- OBI or BAO?

# FDA GSRS Use Case

Which disease ontology to use for annotating clinicaltrial.gov data?

By Alex Welsh and Larry Callahan (FDA)

# FDA's Global Substance Registration System

*Home of the Unique Ingredient Identifier (UNII)*



- **Substance**: Any matter of defined composition that has discrete existence, whose origin may be biological, mineral or chemical. (ISO 11238)
- International collaborative.
- NCATS: GINAS
- ~300K substance to clinical trial relationships via "intervention" in clinicaltrials.gov

## Quick Links

### Substances ^

👓 Browse Substances

🔍 Structure Search

🔍 Sequence Search

🔍 Advanced Search

### Other ⌄



# Global Substance Registration System - GSRS

The main goal of the GSRS software is to assist agencies in registering and documenting information about substances found in medicines. The Global Ingredient Archival System provides a common identifier for all of the substances used in medicinal products, utilizing a consistent definition of substances globally, including active substances under clinical investigation, consistent with the ISO 11238 standard.

Search Substances 🔍

# Mapping to Clinical Trial enables links from application to trials

**Substance Hierarchy**

🔍 **GEFITINIB**                                   S65743JHBS
                                                 *{ACTIVE MOIETY}*

| Application Count: | Product Count: | Clinical Trial Count: | Adverse Event Count: |
|---|---|---|---|
| CDER GSRS: 134 | Active: 18 | 329 | 7659 |
| | Inactive: 0 | | |

**Substance Hierarchy**

▷ 🔍 **NAZARTINIB**                                KE7K32EME8
                                                 *{ACTIVE MOIETY}*

| Application Count: | Product Count: | Clinical Trial Count: | Adverse Event Count: |
|---|---|---|---|
| CDER GSRS: 5 | Active: 0 | 9 | 17 |
| | Inactive: 0 | | |

# Linking conditions to substance and trials is challenging!

# Possible solutions

- Make use of ClinicalTrials.gov strategies to categorize trials by condition, inside the GSRS.

- Use NLP strategies to classify raw clinical trial conditions text into a **broad** set of organ::disease terms for easier searching and faceting.

ClinicalTrials.gov

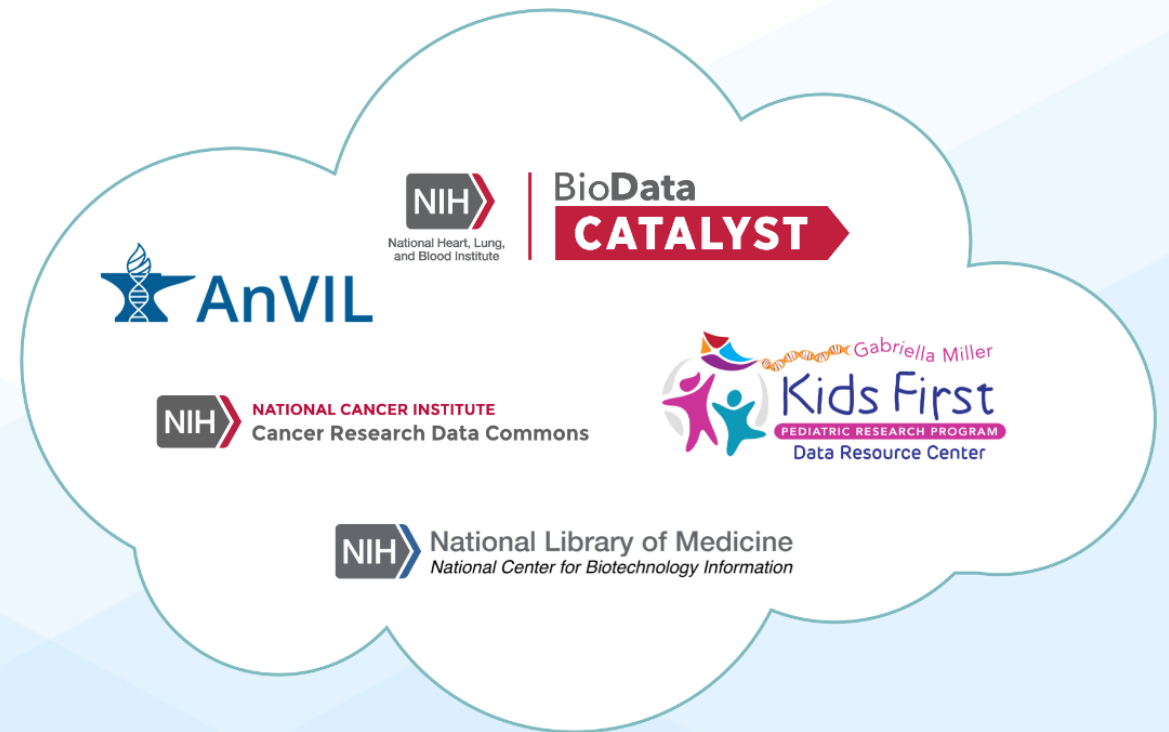| Terms | Search Results* | Entire Database** |
|---|---|---|
| Synonyms | | |
| **Pulmonary Neoplasm** | 8,499 studies | 8,499 studies |
| Lung Cancer | 7,755 studies | 7,755 studies |
| Lung Neoplasm | 7,081 studies | 7,081 studies |
| Lung carcinoma | 1,270 studies | 1,270 studies |
| lung tumors | 90 studies | 90 studies |
| Carcinoma of the Lung | 48 studies | 48 studies |
| Cancer of the Lung | 41 studies | 41 studies |
| Neoplasm of lung | 35 studies | 35 studies |
| CARCINOMA OF LUNG | 26 studies | 26 studies |

# NCPI Use Case

Which ontology to use for disease hierarchy in facet search display?

# What is NCPI?

The NIH Cloud Platform Interoperability (NCPI) effort aims to establish and implement guidelines and technical standards to empower end-user analyses across participating NIH cloud platforms, to facilitate the realization of a trans-NIH, federated data ecosystem.

Established in late 2019 as a coalition of independently funded NIH IC cloud-based data platforms, with additional support from ODSS

https://anvilproject.org/ncpi



By Valentina di Francesco (NHGRI)

# Mapping the diseases to MONDO/DO/HPO

- Only look for "exact match" or "equivalent match"
- MESH IDs to DO IDs using DO-MESH mapping file provided by DO(Lynn Schriml), then manually evaluated using manual search on DO.
- MONDO mapping: using DO mapped IDs to find MONDO IDs, then search the disease names for MONDO IDs.
- If neither MONDO or DO, then search HPO. If not HPO, then search NCIT.

# Mapping Results:

- **73% (43/59)** is mapped to DO IDs
- **78% (46/59)** mapped to MONDO IDs (includes all 43 DO IDs)
- HPO only term (2):
  Venous thrombosis, Left ventricular hypertrophy
- NCIT only term (1):
  Prostatic Neoplasms, Castration-Resistant
- 17% (10/59) terms are not mapped (**non disease terms**):
  Arterial Pressure, Blood Pressure, Lipids, Mendelian Conditions, Metabolomics, Platelet Aggregation, Population, Reference Values, Women's Health, Xenograft Model Antitumor Assays

# Decisions for the NCPI dataset catalog

- Remain using the MeSH terms.

- Plan to use MeSH hierarchy.

- Non disease terms can not be covered by any of the candidate ontologies.

- Switch to display as "disease/focus" as the same in dbGaP.

- Arteriosclerosis, Coronary

Previous Indexing:

- Coronary Disease (1966-1986)

See Also:

- Atherectomy, Coronary

All MeSH Categories
  Diseases Category
    Cardiovascular Diseases
      Heart Diseases
        Myocardial Ischemia
          Coronary Disease
            **Coronary Artery Disease**

All MeSH Categories
  Diseases Category
    Cardiovascular Diseases
      Vascular Diseases
        Arterial Occlusive Diseases
          Arteriosclerosis
            **Coronary Artery Disease**

All MeSH Categories
  Diseases Category
    Cardiovascular Diseases
      Vascular Diseases
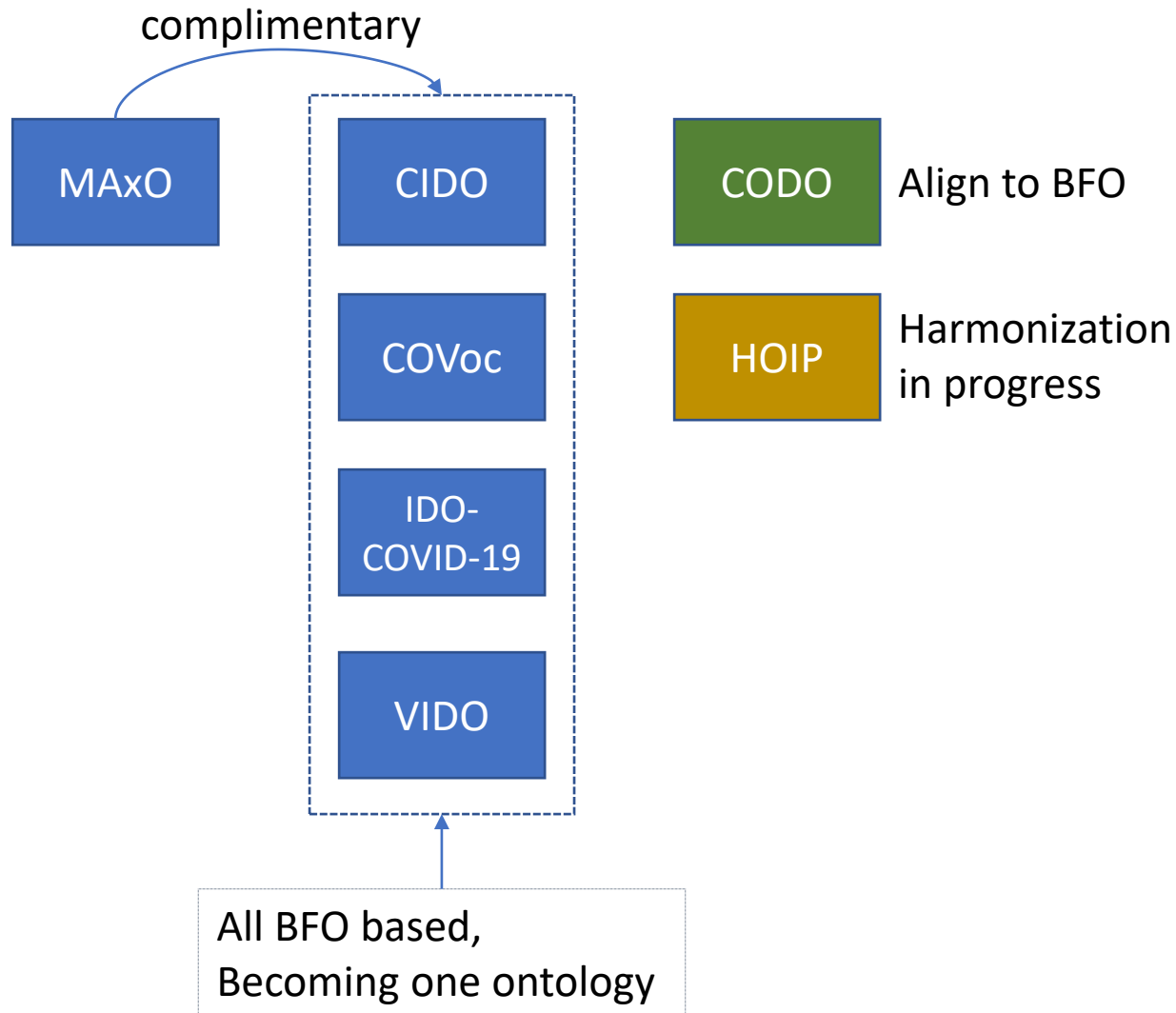        Myocardial Ischemia
          Coronary Disease
            **Coronary Artery Disease**

# What will end users need?

- A single system encompass MONDO, DO, HPO, MeSH, and NCIT
- Mapping to other standard systems, such as OMOP vocabulary, SNOMED CT, MedDRA, etc.
- Reliable "exact match" and "equivalent match" synced in all relevant ontologies.
-  Weighted matches for user's references.
- A unified translator mid-layer to point to this single system.
- Channels to feedback to ontology developers, and more importantly, **knowing** that end users can provide feedback to ontology community and submit new terms.

# Possible solutions

- Ontology harmonization: COVID-19 ontology harmonization as an example.
  - Pro: less ontology mapping, one term for a concept
  - Cons: time consuming, requires higher coordination of ontology developers, may not be realistic
- Ontology mappings: standardized metrics?
- Tools support : enhance end users and developer's interactions.
- Unified interface to interact with end users, build a business service model to respond to end user's request, educational material to lower the entrance bar for users.

# Thank You!

asiyah.lin@nih.gov