Nov. 9. Daegu, South Korea

Creating a Biomedical Data Standardization Ecosystem: Ontologies, Metadata, and Applications for Effective Data Management

Xiaolin Yang

yangxl@pumc.edu.cn

Chinese Academy of Medical Sciences & Peking Union Medical College
BioMedical Branch of China National Population Health Data Center (BMICC)
2023-11-9





Outline



- Goals and status quo of scientific data management
- Biomedical Data Standards Environment
- Standard resource construction
- Summary



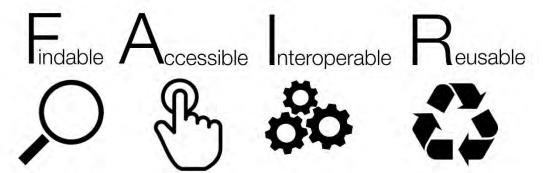
SCIENTIFIC DATA 1101101

SUBJECT CATEGORIES

» Publication characteristics

OPEN Comment: The FAIR Guiding » Research data Principles for scientific data management and stewardship

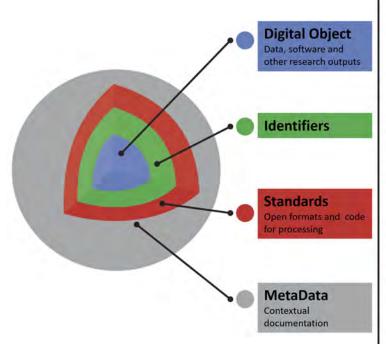
Mark D. Wilkinson et al.#



FAIR Principles put specific emphasis on enhancing the ability of machines to automatically find and use the data ...



Details of the FAIR Principles



Box 2 | The FAIR Guiding Principles

To he Findable:

- F1. (meta) data are assigned a globally unique and persistent identifier
- F2. data are described with rich metadata (defined by R1 below)
- F3. metadata :learly and explicitly include the identifier of the data it describes
- F4. (meta)data are registered or indexed in a searchable resource

To be Accessible:

- A1. (meta) data are retrievable by their identifier using a standardized communications protocol
- A1.1 the protocol is open, free, and universally implementable
- A1.2 the protocol allows for an authentication and authorization procedure, where necessary
- A2. metadata are accessible, even when the data are no longer available

To be Interoperable:

- 11. (meta data use a formal, accessible, shared, and broadly applicable language for knowledge representation
- 12. (meta)data use vocabularies that follow FAIR principles
- 13. (meta) pata include qualified references to other (meta) data

To be Reusable:

- R1. meta data) are richly described with a plurality of accurate and relevant attributes
- R1.1. (meta)data are released with a clear and accessible data usage license
- R1.2. (meta)data are associated with detailed provenance
- R1.3. (meta)data meet domain-relevant community standards



Scientific data Content standards



Conceptual model, schema, exchange formats etc

- Define the structure and interrelation of information, and the transmission format
- e.g. FASTA, VCF





Minimum information reporting requirements, checklists

- Report the same core, essential information
- o e.g. MIAME guidelines

Controlled vocabularies, taxonomies, thesauri, ontologies etc.

- Provide definitions and unambiguous identification or concepts and objects.
- e.g. Gene Ontology



Formal systems for resources and other digital objects that allow their unique and unambiguous identification.

DOI, ORCID iD Identifier Schema,
QID (Wikidata Identifier)







Minimum Information About a Microarray Experiment - MIAME

MIAME describes the Minimum Information About a Microarray Experiment that is needed to enable the interpretation of the results of the experiment unambiguously and potentially to reproduce the experiment. [Brazma et al., Nature Genetics]

The six most critical elements contributing towards MIAME are:

- 1. The raw data for each hybridisation (e.g., CEL or GPR files)
- The final processed (normalised) data for the set of hybridisations in the experiment (study) (e.g., the gene expression data matrix used to draw the conclusions from the study)
- The essential sample annotation including experimental factors and their values (e.g., compound and dose in a dose response experiment)
- The experimental design including sample data relationships (e.g., which raw data file relates to which sample, which hybridisations are technical, which are biological replicates)
- Sufficient annotation of the array (e.g., gene identifiers, genomic coordinates, probe oligonucleotide sequences or reference commercial array catalog number)
- The essential laboratory and data processing protocols (e.g., what normalisation method has been used to obtain the final processed data)

For more details, see MIAME 2.0.

MIAME does not specify a particular format, however, obviously the data are more usable, if it is encoded in a way that the essential information specified by MIAME can be accessed easily. FGED recommends the use of MAGE-TAB format, which is based on spreadsheets, or MAGE-ML.

MIAME also does not specify any particular terminology, however for automated data exchange the use of standard controlled vocabularies and ontologies are desirable. FGED recommends the use of MGED Ontology for the description of the key experimental concepts, and where possible ontologies developed by the respective community for describing terms such as anatomy, disease, chemical compounds etc (see OBO page for more detail).













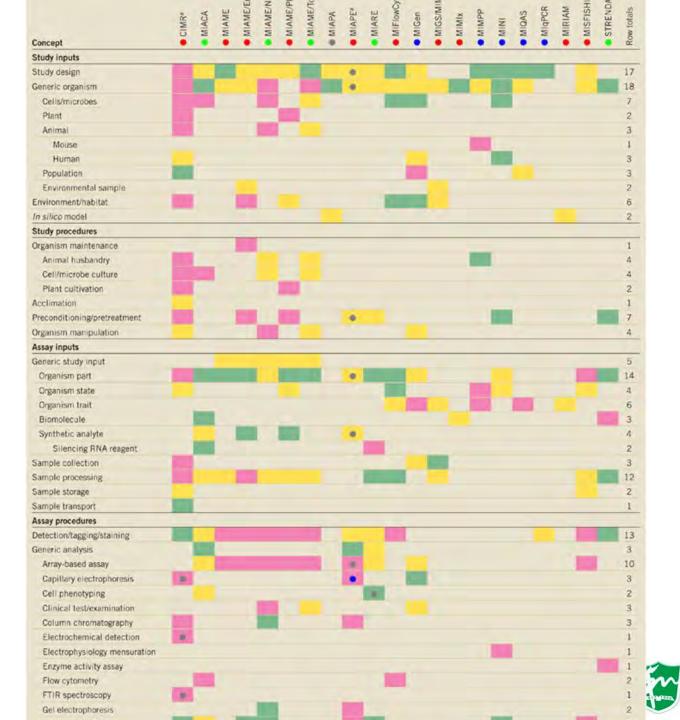




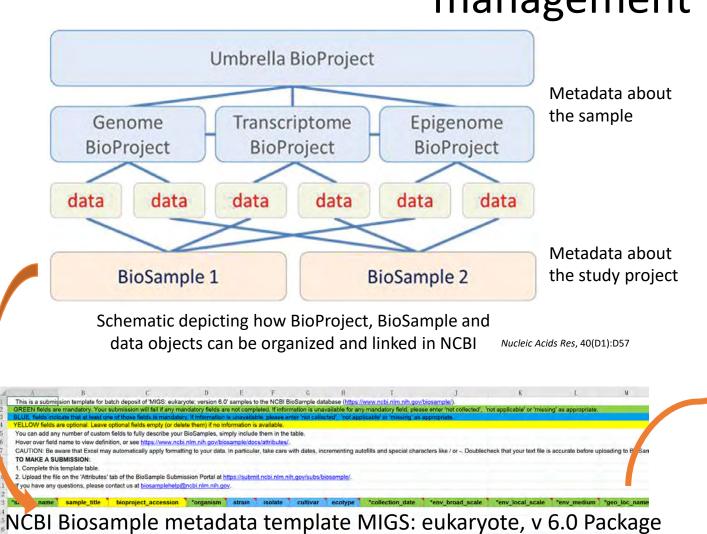


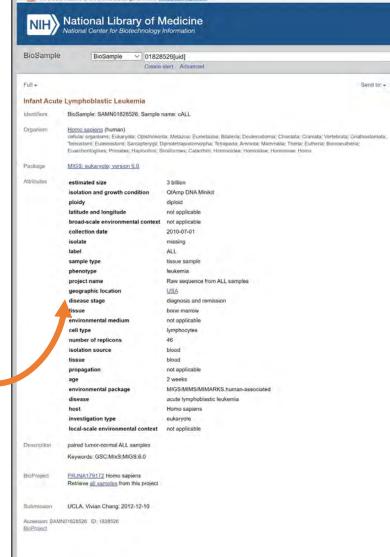






Example of NCBI primary database metadata management







Metadata Quality in NCBI BioSample

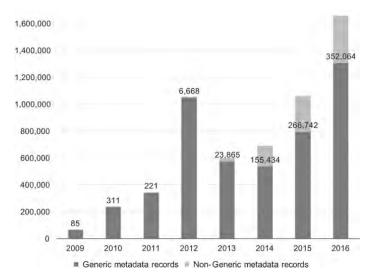


Figure 3. Metadata submissions to NCBI BioSample from 2009-2017.

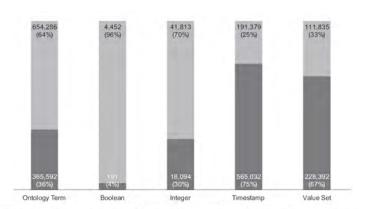


Figure 5. Quality of attributes in packaged metadata records in NCBI BioSample. The columns represent the metadata attribute types. Each column shows the number and percentage of metadata attributes whose values are either well-specified or invalid.

- Majority of data submitters use generic metadata templates, which lack adequate and accurate descriptions of sample attributes
 - Among the data elements (18,650) that characterize sample attributes in the BioSample database, only 2.4% (452) of the data elements (sample attributes) specified by NCBI and 97.6% are user-defined data elements
 - Only 9 of them are named using controled terminology
- Of the sample attribute values that require the use of ontology, only 36% use ontology terms in at least one value

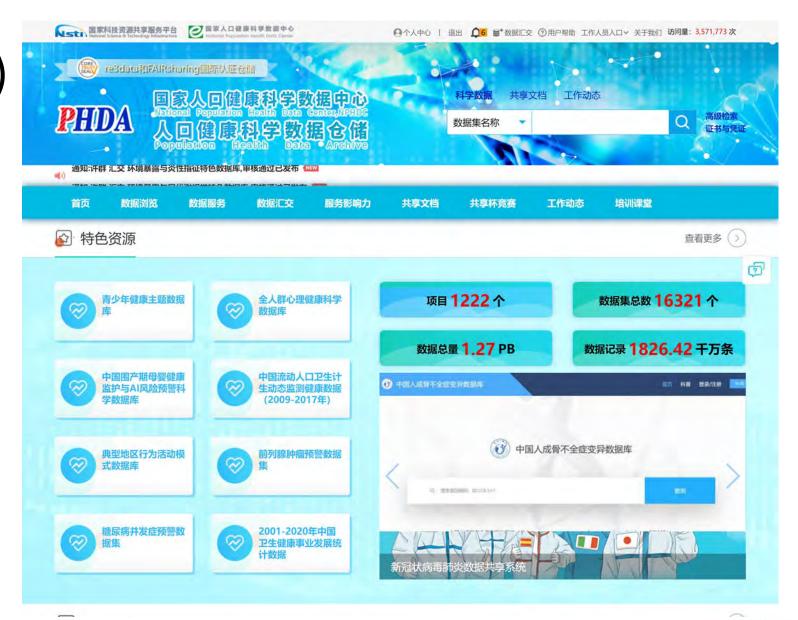
Progress in Scientific Data Management and Sharing in China

- *Measures for the Management of Scientific Data* were issued by the Chinese government in 2018
- Scientific data generated from government-funded science and technology projects should be submitted to the designated data center.
- National Population Health Data Center (NPHDC, formerly named NCMI) is one of the nationally recognized data centers in biomedical and population health field
 - Main responsibility is receiving and archiving data that scientists submit
 - Developing domain data standards
 - Providing data sharing and reuse services



Population Health Data Archive (PHDA)

- https://www.ncmi.cn/
- The new system was available in 2019.
- PHDA is an infrastructure to ensure domestic researcher share their data.
- Provide a platform for querying and reusing the data

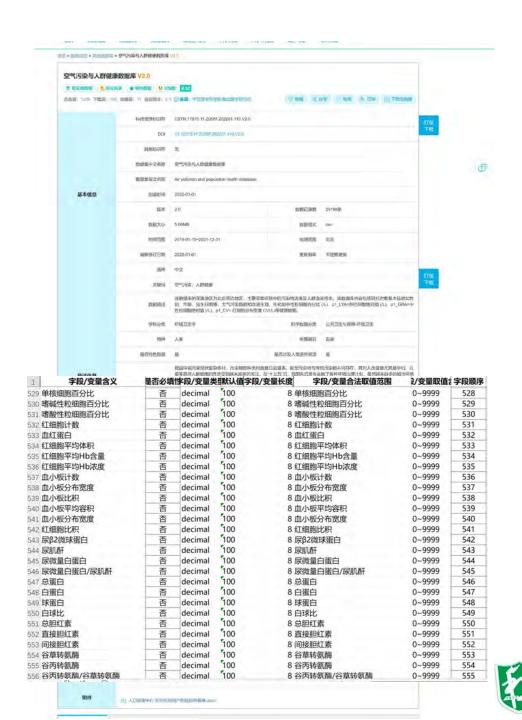




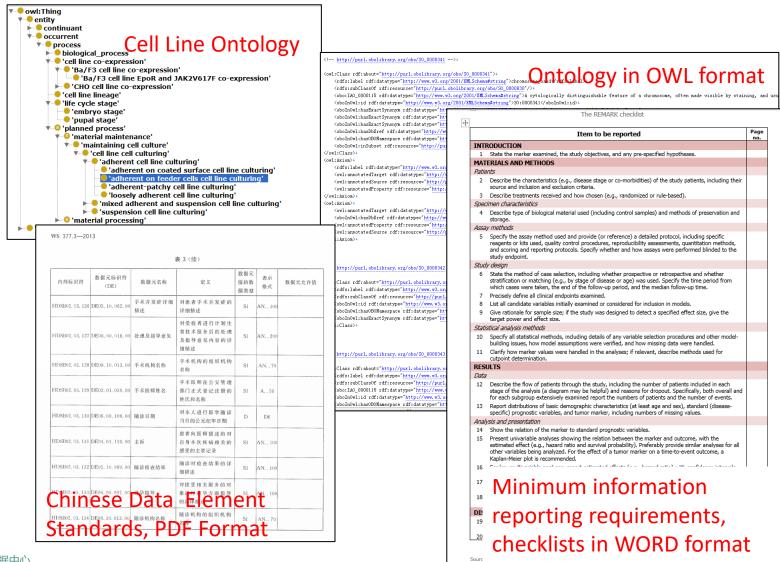
Challenges for PHDA

- The submitted data with complicated types from multiple sub-domains.
- Project-centric data submitting is useful for aggregating data in a short time, we still need to promote data in-depth reuse.
- Metadata submitted as data dictionary without pre-defined rules.
 - Without any semantics standard
- Difficulties in harmonizing data according to specific needs.





Data standards are difficult for domain scientists to use



- Hard to find
- Hard to obtain
- Hard to understand
- Hard to use



Outline

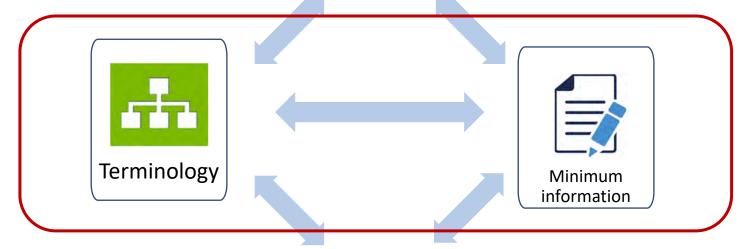
 Goals and status quo of scientific data management



- Biomedical Data Standards Environment
- Standard resource construction
- Summary



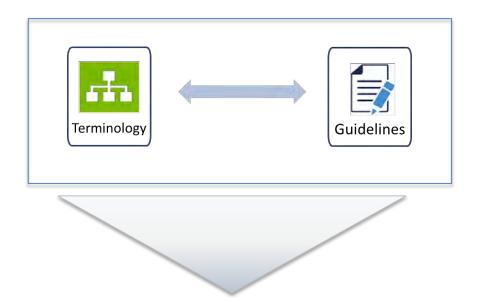




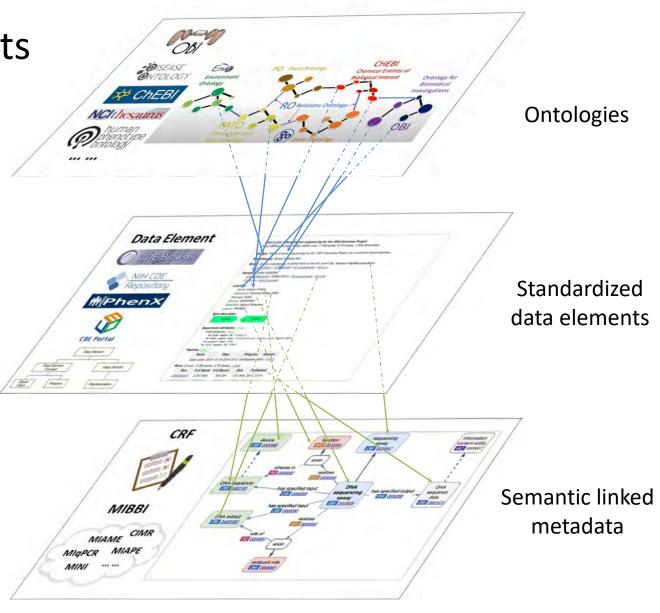




From ontology and data elements to metadata standardization



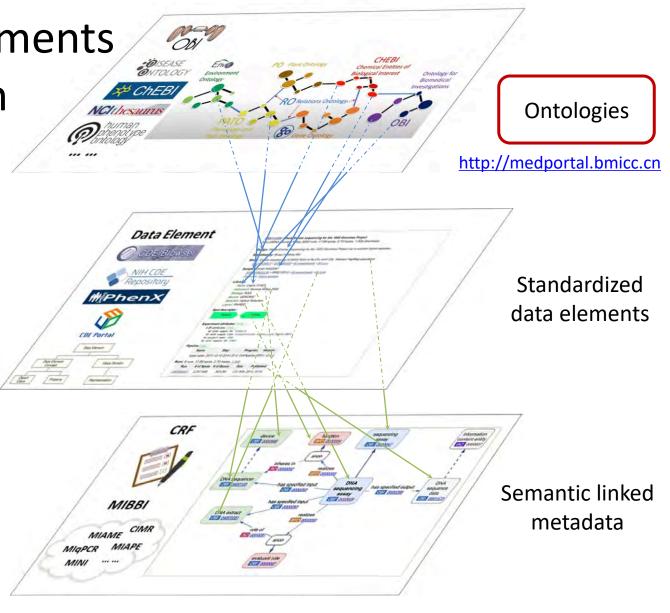
Machine-understandable content metadata standards



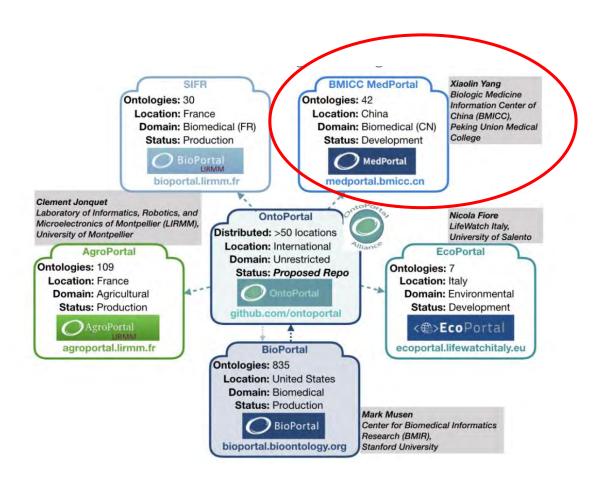
From ontology and data elements to metadata standardization

Terminology

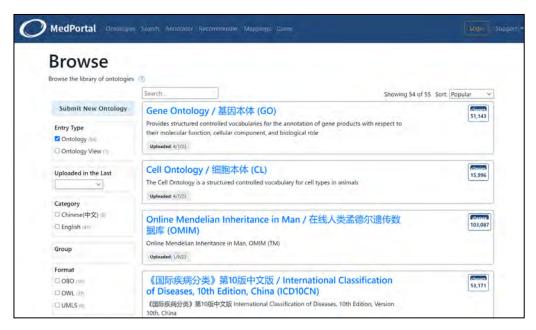
Machine-understandable content metadata standards



MedPortal – a repository of biomedical ontology resources

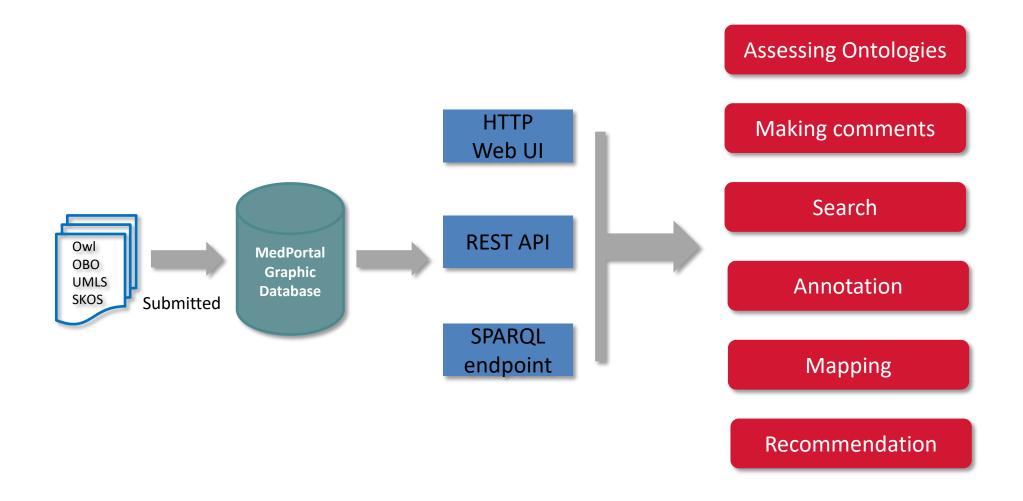


- Reuse the Ontoportal Alliance (BioPortal) framework
- Provide bilingual ontology services in Chinese and English
- Ontology services that support semantic standards

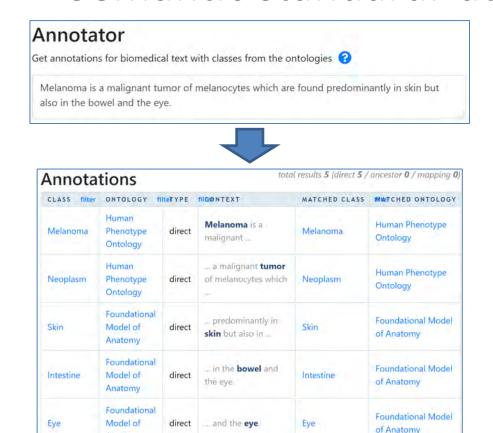


URL: http://medportal.bmicc.cn/

Basic Functions of MedPortal



Semantic standardization with MedPortal





REST API

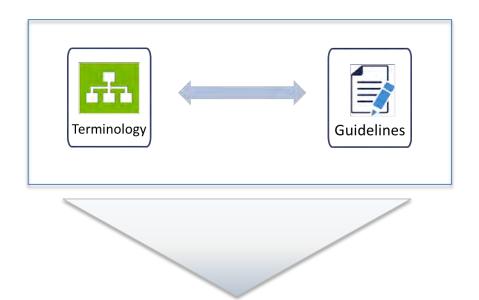
- Annotate text
 - o GET POST /annotator?text={input text}
 - example: /annotator?

Anatomy

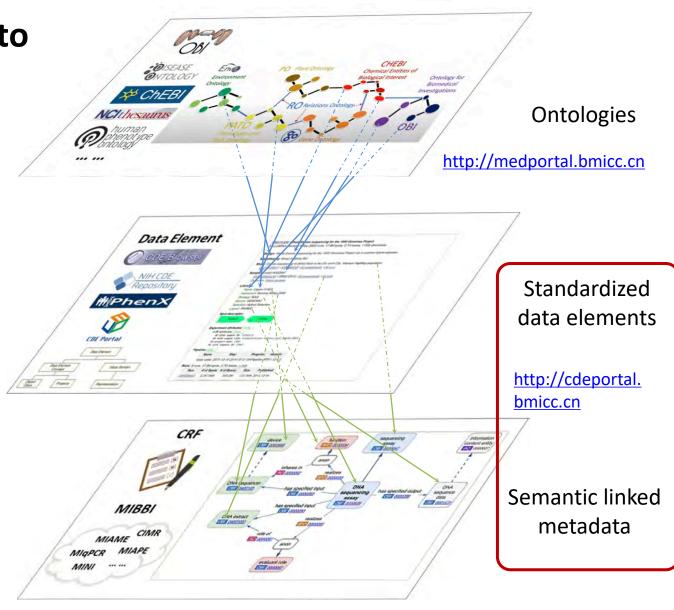
– Previous 123456789...3132Next → HUMAN PHENOTYPE ONTOLOGY EXPERIMENTAL FACTOR ONTOLOGY SOURCE Corneal neovascularization LOOM corneal neovascularization Prelingual sensorineural hearing impairment Prelingual sensorineural hearing impairment SAME URI Hyperkalemia Hyperkalemia SAME URI adrenocortical adenoma Adrenocortical adenoma LOOM Decreased liver function Decreased liver function SAME_URI SAME_URI Breech presentation Breech presentation anxiety Anxiety LOOM Oral ulcer Oral ulcer SAME_URI Diffuse mesangial sclerosis Diffuse mesangial sclerosis SAME_URI SAME URI Talipes equinovarus Talipes equinovarus Cyclic neutropenia Cyclic neutropenia LOOM papillary thyroid carcinoma Papillary thyroid carcinoma LOOM Breast mass SAME_URI MOON Cunarficial enreading malanoms

Ontology terminology mapping

From ontology and data elements to metadata standardization

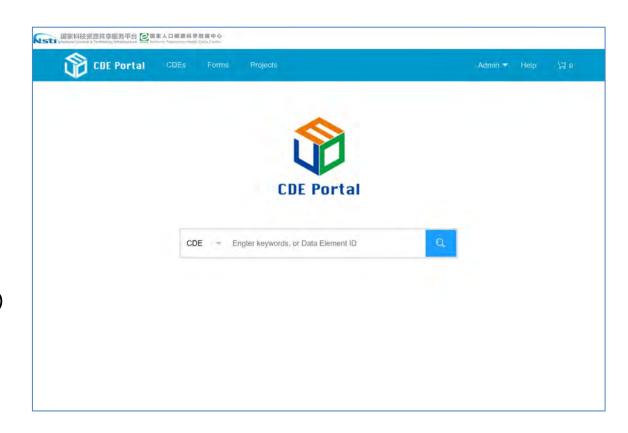


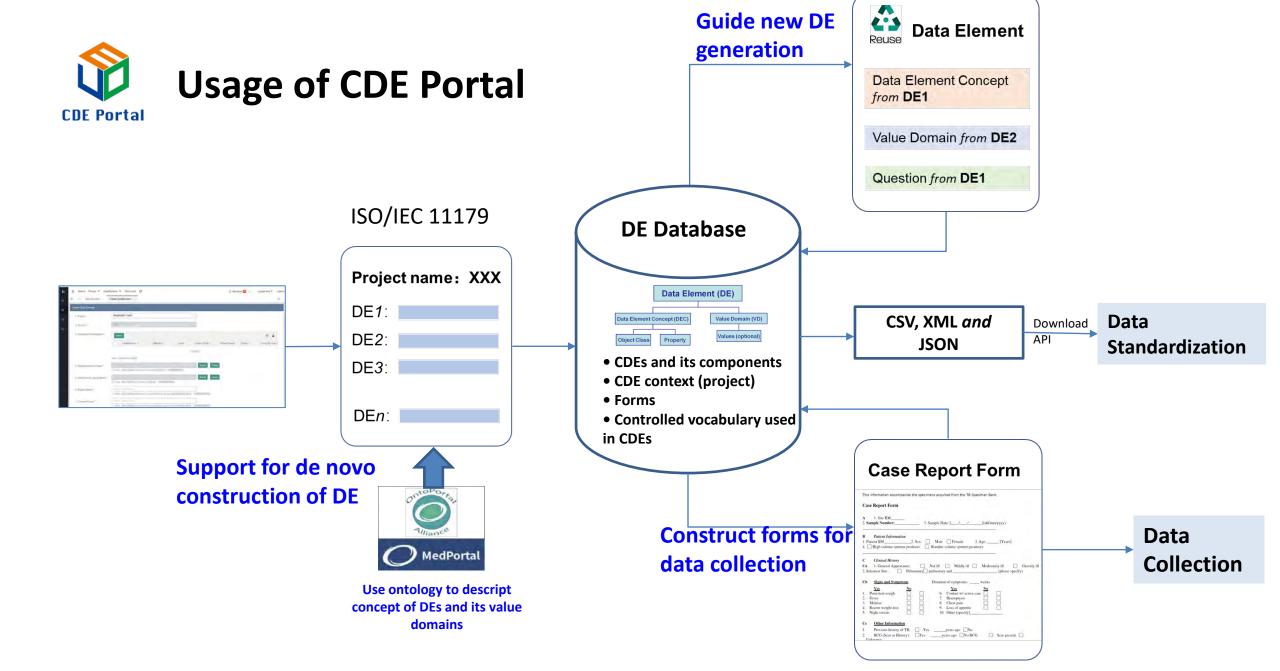
Machine-understandable content metadata standards



CDE Portal

- https://cdeportal.bmicc.cn/
- Provide tools for metadata compilation, registration, review, and publication for the NCMI
 - Tools for data element building
 - Tools for reusing data elements
- Provide high-quality, common data elements (Common Data Element, CDE)
 - Data elements for important research projects at home and abroad
 - Important, standardized domain data elements





Represent Data Elements in a FAIR way



 Cancer Generic Prima 	ry Tumor TNM Finding/癌症	定一般性原发肿瘤TNM分期 Standard
A cancer finding in the TNM sy	stem that is relevant to the diag	nosis of cancer.
Preferred Question:Cancer (Generic Primary Tumor TNM Fin	ding
Value Domain:Generic Prima	ry Tumor TNM Finding	
	G : D: T TMM	Fledber
Data Element Concept:Canc	er Generic Primary Tumor TNM	Finding
		F
		ow: Draft new Create Date: 2021-10-16
		TO THE ALL ADVISOR OF
Public ID: DE0071220 Version	on: 1.0 Project: BMICC workflo	TO THE ALL ADVISOR OF
Public ID: DE0071220 /ersic	on: 1.0 Project: BMICC workflo	TO THE ALL ADVISOR OF
Public ID: DE0071220 /ersicopermissible Value Any T	on: 1.0 Project: BMICC workflow Value Meaning Any T	TO THE ALL ADVISOR OF

Implement semantic support for the definition and value of data elements Object Class Concepts Chinese **English Name** Definition Concept Code ontology Primary Name A cancer finding in the TNM system that is relevant to Cancer TNM http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#C48232 None Finding the diagnosis of cancer. Thesaurus Property Concepts Definition Concept Code Ontology Primary **English Name** Chinese Name Generic Primary Tumor TNM Finding http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#C48885 **NCI Thesaurus** Representation Concepts **English Name** Chinese Name Definition Concept Code Primary Ontology Generic Primary Tumor TNM Finding None http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#C48885 NCI Thesaurus Permissible Values Begin Date -Public ID Value Project Meaning End Date PV0102942 **BMICC** Any T No specification on the details of the primary tumor growth 10-16 -T0 Stage 2021-NCI PV0035492 A primary tumor TNM finding indicating that there is no evidence of primary tumor. Finding caDSR T1 Stage 2021-NCI PV0035493 A clinical and/or pathologic primary tumor TNM finding indicating that the cancer is limited to the site of growth. Finding 10-16 caDSR A general term that refers to a TNM finding of a primary tumor limited to the site of growth. The definition of T1a TNM T1a finding depends on the specific type of cancer that it refers to; for example, for breast cancer it refers to a primary tumor 2021-NCI PV0035494 Stage that is more than 0.1cm, but not more than 0.5 cm in greatest dimension; for kidney cancer it refers to a primary tumor caDSR Finding that is 4 cm or less in greatest dimension; and for thyroid cancer it refers to a primary tumor that is 1 cm or less in

A general term that refers to a TNM finding of a primary tumor limited to the site of growth. The definition of T1b TNM

finding depends on the specific type of cancer that it refers to; for example, for breast cancer it refers to a primary tumor

that is more than 0.5 cm, but not more than 1.0 cm in greatest dimension; for kidney cancer it refers to a primary tumor

that is more than 4 cm, but not more than 7 cm in greatest dimension; and for thyroid cancer it refers to a primary tumor

greatest dimension.

T₁b

Stage

Finding

T1b

PV0035495

urus.owl#C48885

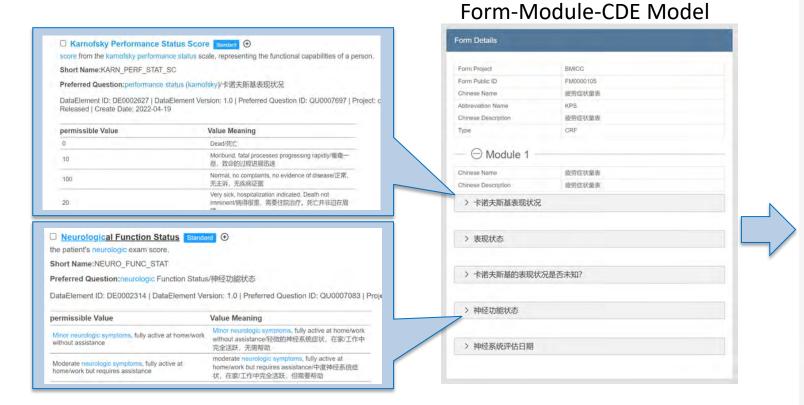
24

2021-

NCI

caDSR

Case investigation forms (CRFs) facilitate the reuse of data elements and improve data interoperability



Provide precise and interoperable metadata and form for re-using

API

CDE and Form can be reused in JSON/CSV

皮劳症状量表

疲劳症状量表

疲劳症状量表

	死亡	
	奄奄一息, 致命的过程进展迅速	
	正常,无主诉,无疾病证据	
	病得很重,需要住院治疗。死亡并非迫在眉睫	
	严重残疾,需要住院治疗。死亡并非迫在眉睫	
	残疾,需要特别照顾和帮助	
	需要大量援助和频繁的医疗护理	
	需要偶尔的帮助,但能够照顾他她的大部分需求	
	照顾自己, 无法进行正常活动或积极工作	
	正常活动与努力;疾病的一些体征或症状	
	能够进行正常活动;轻微的疾病体征或症状	
	现状态 送土斯其的主现状况显示主知?	
	现状态 诺夫斯基的表现状况是否未知? 否	
	诺夫斯基的表现状况是否未知?	
3. 卡	诺夫斯基的表现状况是否未知? 否	
3. 卡	诺夫斯基的表现状况是否未知? 否 是	
3. 卡	诺夫斯基的表现状况是否未知? 否 是 经功能状态	
3. 卡	诺夫斯基的表现状况是否未知? 否 是 经功能状态 轻微的神经系统症状,在家/工作中完全活跃,无需帮助	
3. 卡	诺夫斯基的表现状况是否未知? 否 是 经功能状态 轻微的神经系统症状,在家/工作中完全活跃,无需帮助中度神经系统症状,在家/工作中完全活跃,但需要帮助	
3. 卡	诺夫斯基的表现状况是否未知? 否 是 经功能状态 轻微的神经系统症状,在家/工作中完全活跃,无需帮助 中度神经系统症状,在家/工作中完全活跃,但需要帮助 中度神经系统症状,在家/工作中缺乏完全活动,需要安抚	后法
3. 卡4. 神	诺夫斯基的表现状况是否未知? 否 是 经功能状态 轻微的神经系统症状,在家/工作中完全活跃,无需帮助中度神经系统症状,在家/工作中完全活跃,但需要帮助中度神经系统症状,在家/工作中缺乏完全活动,需要安抚 无神经系统症状,在家/工作中完全活跃,无需帮助	后法

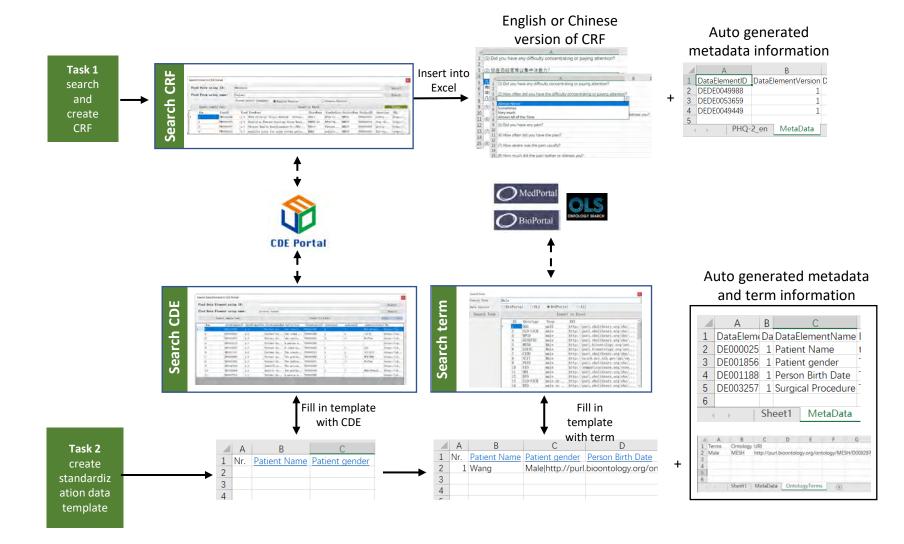
Chinese-English Form

CDE Tools: EXCEL plugin for data standard application

■ URL:

https://github.com/MedportalProject/CDE-Tools

- **■** Features:
 - Excel plug-in
 - Ontology term search
 - Text annotations
 - Batch processing of standardized data elements
- Extensions: Supports all ontology repositories that use OntoAlliance and OLS frameworks



Outline

- Goals and status quo of scientific data management
- Biomedical Data Standards Environment



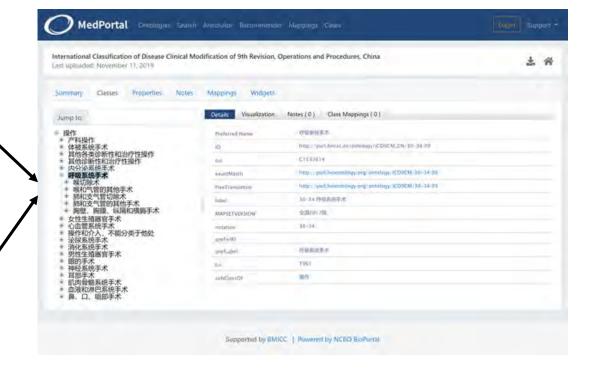
- Standard resource construction
- Summary

Chinese ontology resources in MedPortal

Ontology	Label	Definition	Other Annotation Properties
OGMS	76	70	4
PATO	1742	1722	906
IAO	194	183	116
OBI	1875	1828	422
BFO	36	17	170
RO	475	0	3
ICO	176	161	0
HP	9998	6659	1647
CLO*	878	330	439

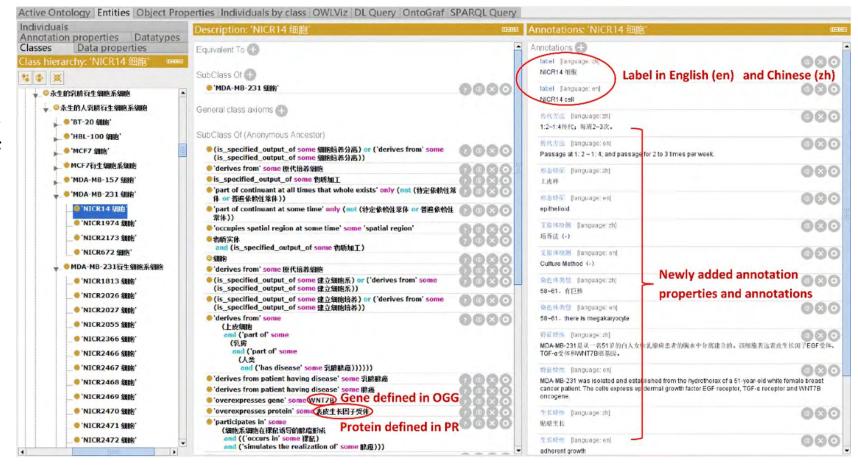
Some important OBO domain ontologies were also translated into Chinese.

Some widely used clinical taxonomies and ontologies, such as ICD-10, ICD-11, LOINC and ICD-9-CM



Cell Line Ontology--Integrate cell line information from the Chinese experimental platform

- Chinese National Infrastructure of Cell Line Resource (NICR, 国家实验细胞资源共享服务平台)
- NICR 2 704 cell lines were integrated
- Added cell line attributes
- Formal representation of multiple biological characteristics of a cell line
- Mapping to cellular resources such as ATCC, Coriell Cell, and JCRB
- bilingual ontology subset



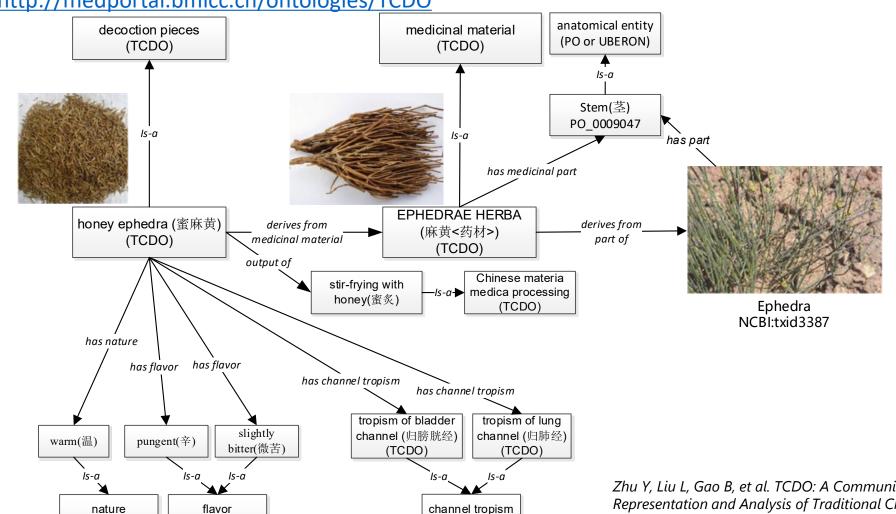
TCDO (Traditional Chinese Drug Ontology)

 Based BFO, more than 400 popular Chinese herbal decoction pieces and Chinese herbal medicines and their attributes are formally represented.

http://medportal.bmicc.cn/ontologies/TCDO

(TCDO)

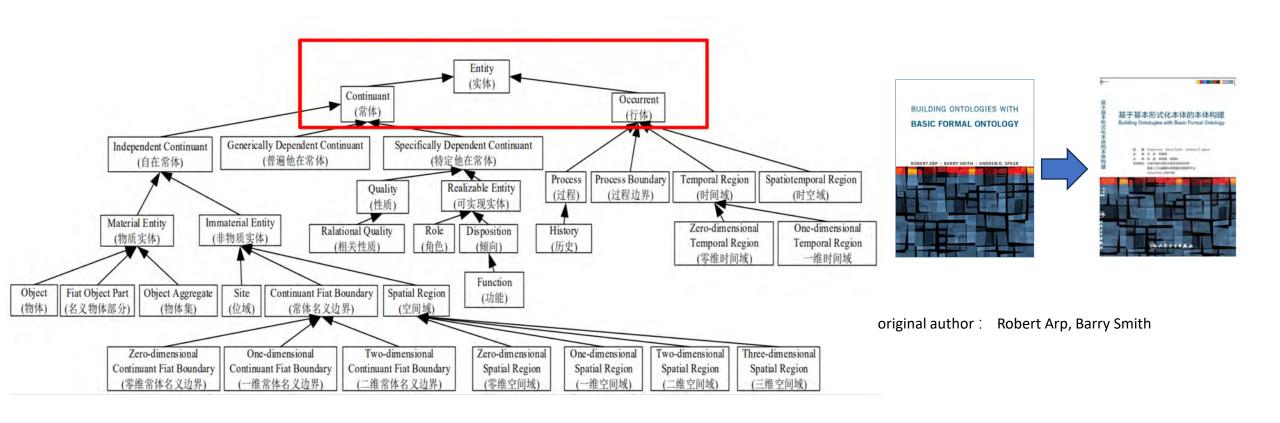
(TCDO)



(TCDO)

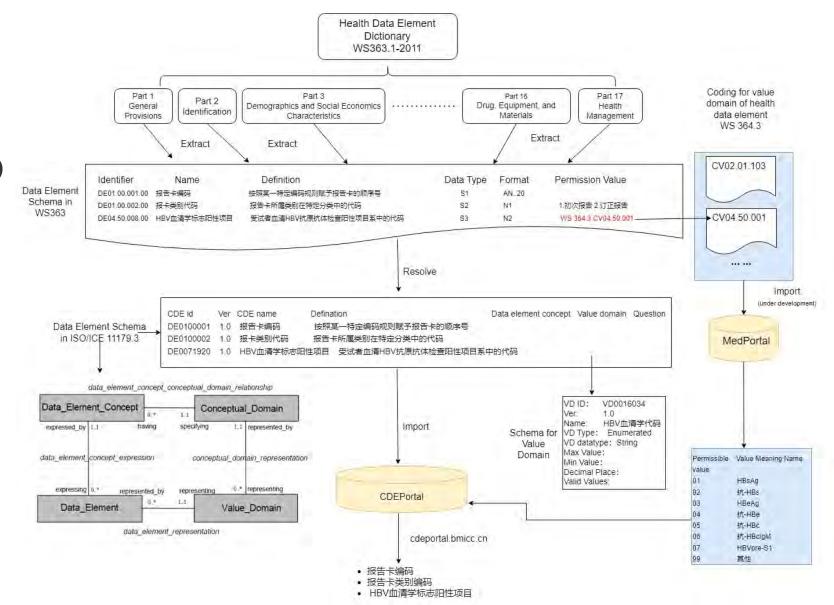
Zhu Y, Liu L, Gao B, et al. TCDO: A Community-Based Ontology for Integrative Representation and Analysis of Traditional Chinese Drugs and Their Properties [J]. Evidence-Based Complementary and Alternative Medicine, 2021, 2021: 6637810.

Translate BFO (Basic Formal Ontology) into Chinese



Health Data Element Dictionary (WS363.1-17&WS 364.1-17)

- Representation of the standards following FAIR guidelines
 - Data elements were resolved and extracted according to ISO11179 and stored in CDE Portal
 - Control vocabulary in the value domain was stored in a value field



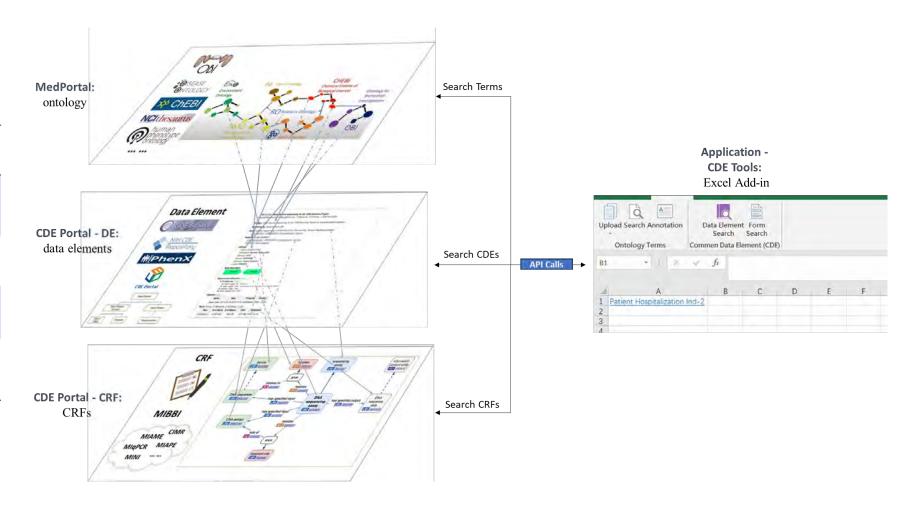
Outline

- Goals and status quo of scientific data management
- Biomedical Data Standards Environment
- Standard resource construction



Summary

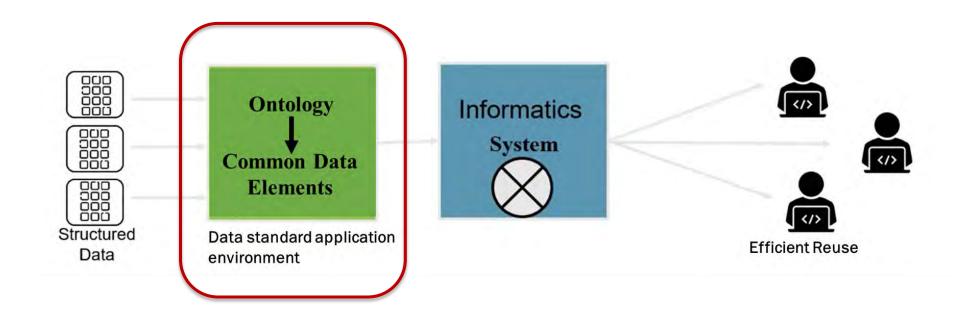
Data Types	Numbers
Ontologies Terms	60 3,583,866
CDEs	72,619
CRFs	106
Value Domain	16,711



Biomedical Data Standards Environment

Next

- Provide support for domain-specific database in PHDA
- Build implementable biomedical data standards by using of MedPortal and CDE Portal



Acknowledgement

- Wang Heng
- Yan Zhu ((IITCM of CACMS)
- Sheng Yang
- Chen Shao
- Hongjie Pan
- Lulu Zhang
- Zhigang Wang
- Zhe Wang
- Yize Yuan
- Jingwen Guo
- Wei Zhou

- Weimin Zhu
- Yongqun Oliver He (Michigan Univ.)
- Jie Zheng(Penns. Univ.)













Thank you!

