



9AM

2PM in London (GMT), 11PM in Tokyo (GMT+9)

Panel: Ontologies and AI

Moderator: Bruce W. Herr II, *Indiana University*

Presenters:

- Maria-Esther Vidal, *German National Library of Science and Technology (TIB), Germany* ([Knowledge Graph-driven hybrid AI](#))
- Yongxin (Kiki) Kong, *Indiana University & Chinese Academy of Sciences, China* ([HRAIit](#))
- Oliver He, *University of Michigan*



Maria-Esther Vidal,
*Leibniz University of Hannover,
TIB-Leibniz Institute of Science and
Technology, University Library
Hannover, Germany*

Knowledge Graph-Driven AI

Maria-Esther Vidal

Professor Data Science Institute, Leibniz University of Hannover

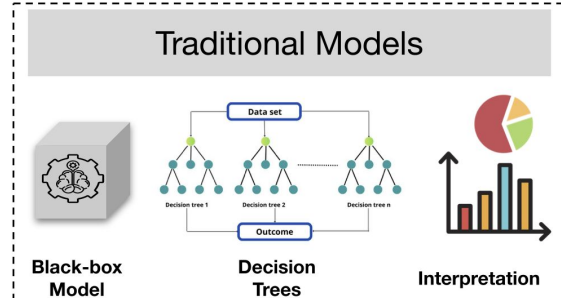
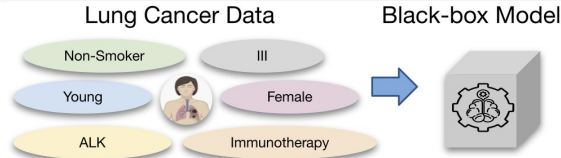
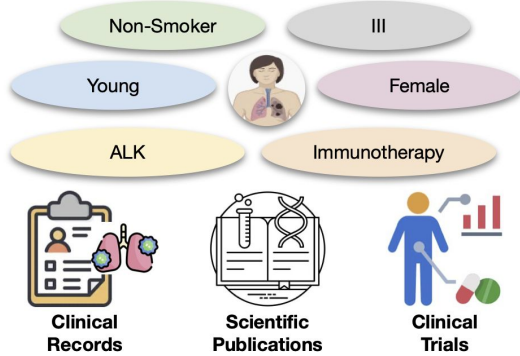
**Head of the Scientific Data Management Group at TIB-Leibniz
Information Center for Science and Technology, Germany**



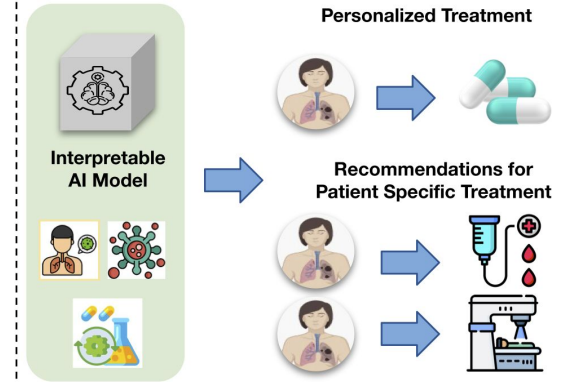
AI Models in Medicine- Scattered Data and Fragmented Knowledge

Negative Impact

Multiple Data Sources and Fragmented Medical Knowledge



Interpretable Models



Clinical Objectives



Clinical objectives:

- Maximize survival time, life quality
- Minimize toxicities, adverse events
- Avoid relapse and disease progression

Decision-making



Decision-making based on:

- data (e.g., symptoms, patient history)
- similar treated patients

Hybrid-AI Systems

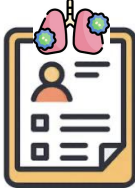


Bridging Cognitive Medical Thinking with AI

- Transparent Decision-Making
- Data Privacy and Sovereignty

Semantic Data Integration- Uniform View of Heterogeneous Data

Data Sources



Lung Cancer Registry



Breast Cancer Registry



Scientific Publications



Scientific Databases



Drug-drug interactions



Genomics and Mutations



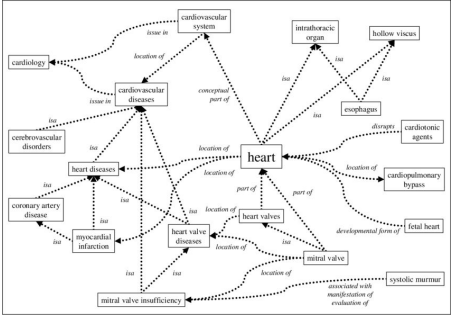
Medical Guidelines and Protocols



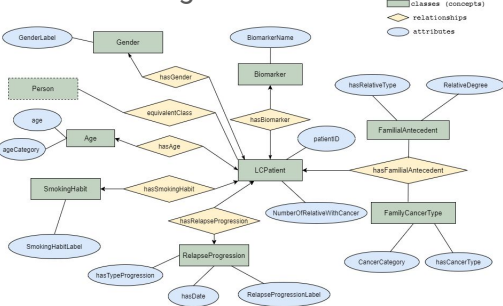
Clinical Trials

Biomedical Knowledge

Unified Medical Language System (UMLS)



Unified ontologies



Data Integration System [1]

A data integration system $DIS = \langle O, S, M, \Sigma \rangle$

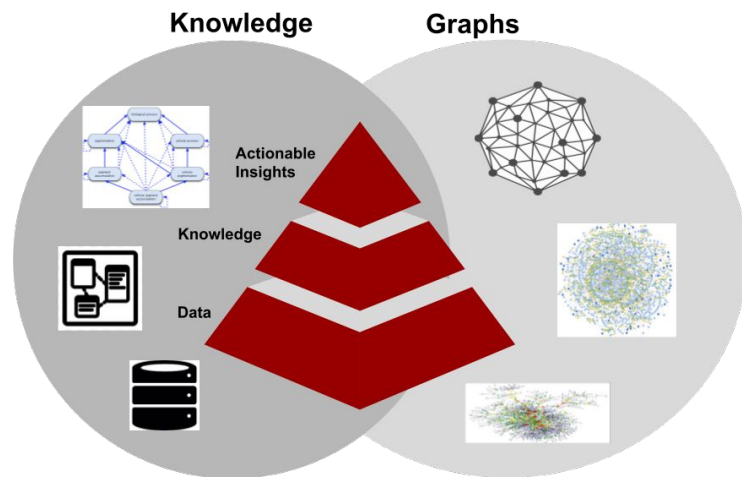
- O is an ontology or schema which provides a uniform view to the data sources in S .
- S is a set of $\{S_1, \dots, S_n\}$ of the signatures of the data sources that compose a DIS.
- M is a set of mappings between signatures of the sources in S and concepts in O .
- $\Sigma = (\varphi, S, \lambda)$ is a shape schema over O ; φ : set of shapes; S : set of shape labels; $\lambda: S \rightarrow \varphi$ total function from labels to shapes.

A data integration system

- provides a uniform view, and to integrate data collected from heterogeneous data sources

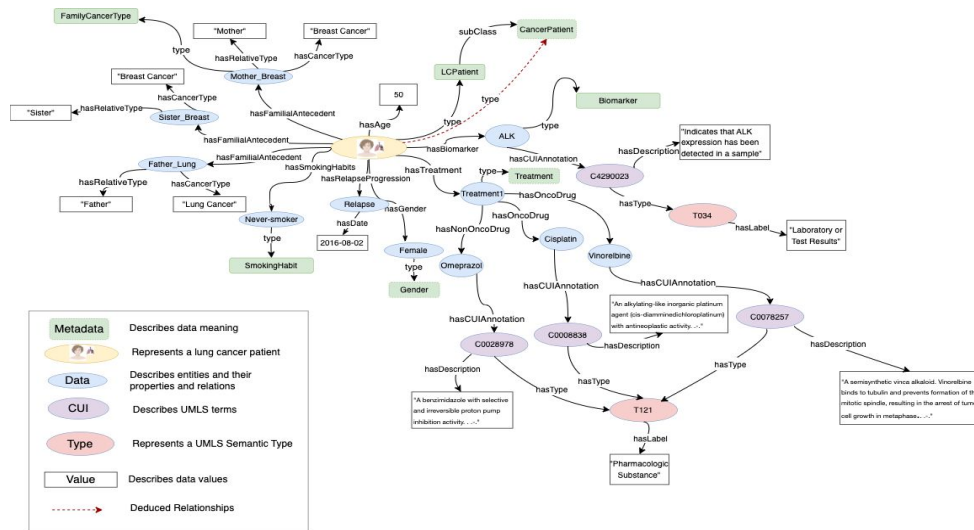
[1] M. Lenzerini. Managing Data through the Lens of an Ontology. AI Mag. 39(2): 65-74 (2018)

Knowledge-Graphs- structures to integrate heterogeneous data, capture domain knowledge, and enable explainable AI through symbolic reasoning



Knowledge Graphs

- data **structures** representing the **convergence** of **knowledge** and **data** as **factual** statements
- using a **graph data model**



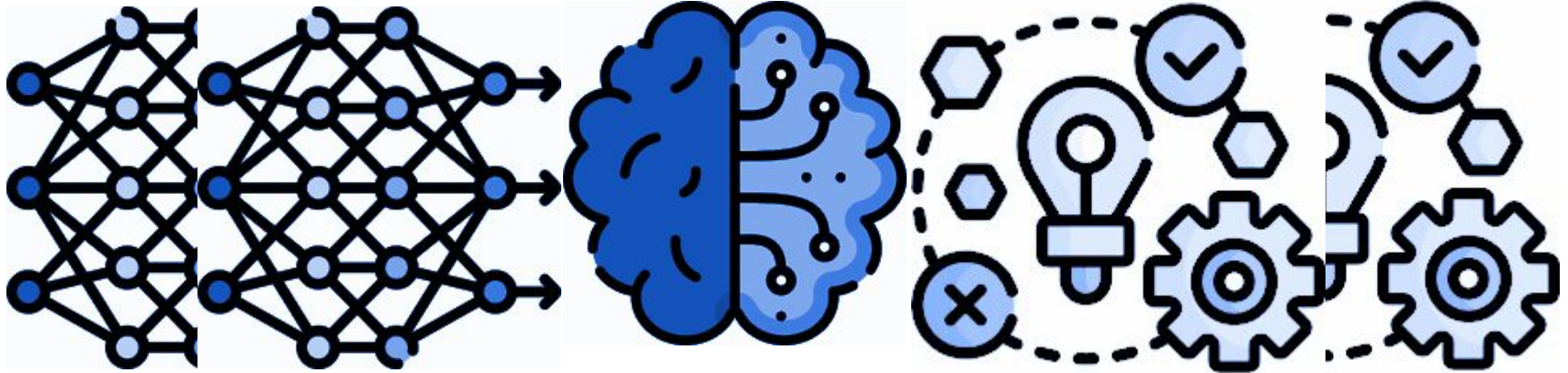
- **Entities** and their **relationships** are represented as **first-class citizens**.
- **Metadata** (via ontologies) describing and providing information about other data.
- **Metadata** and **data** can be **empowered** with **inference** to deduce new facts.

Integrating Semantics and Learning: The Role of Hybrid AI Systems

Neural Components

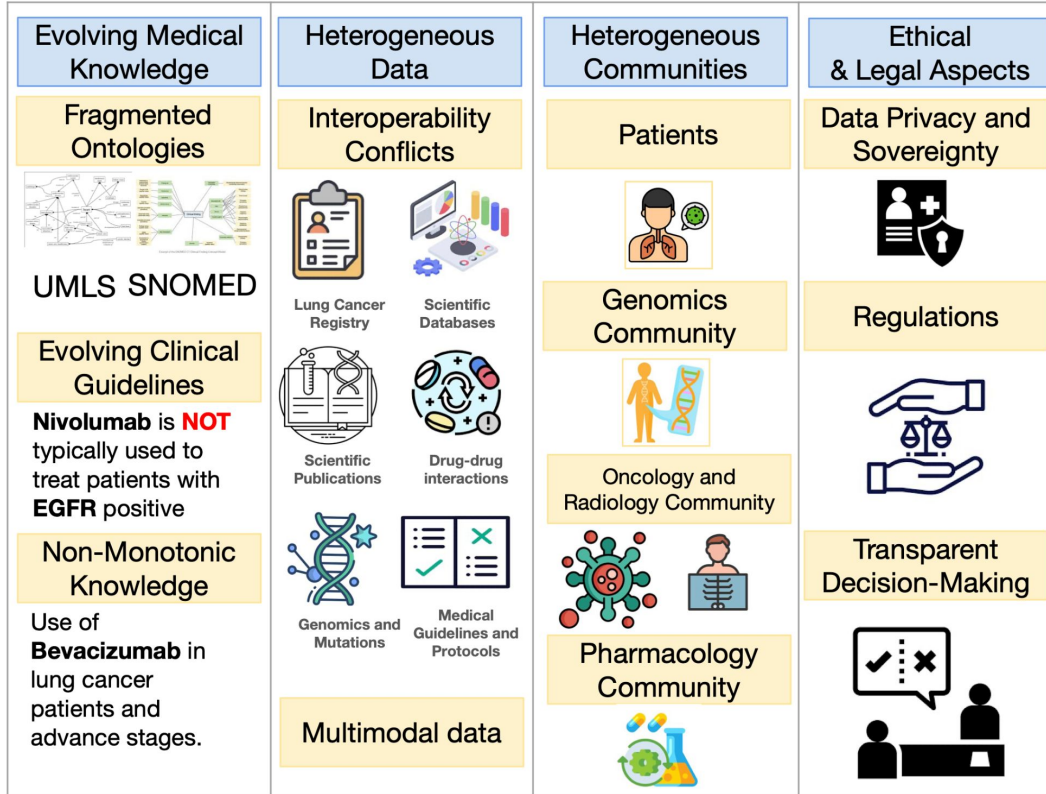
Hybrid AI Systems

Symbolic Components



Data-driven, black-box	AI systems that can learn	logic-based reasoning
that approximate a given	data while also incorporating	modules with explicit
function (e.g., neural networks)	interpretable reasoning	constraints or specifications.

Challenges in Knowledge Graph-driven Hybrid AI Systems



Evolving Medical Knowledge:

- Fragmented ontologies like UMLS and SNOMED limit data consistency.
- Frequent updates to clinical guidelines, e.g., Nivolumab not for EGFR-positive cases.
- Non-monotonic knowledge, e.g., Bevacizumab use in advanced lung cancer.

Heterogeneous Data:

- Multimodal sources: genomic data, drug interactions, and medical guidelines.
- Interoperability issues across registries, databases, and publications.

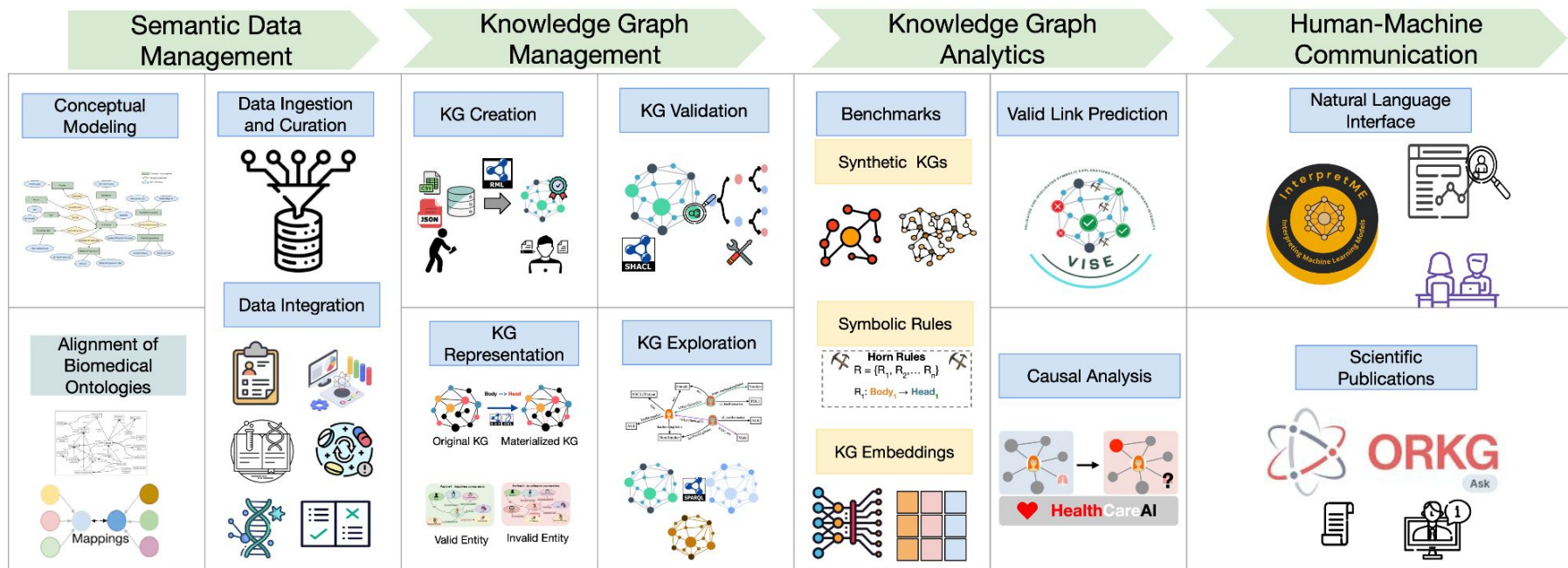
Heterogeneous Communities:

- Stakeholders include patients, genomics experts, oncologists, and pharmacologists.
- Conflicting perspectives complicate shared data modeling.

Ethical and Legal Aspects:

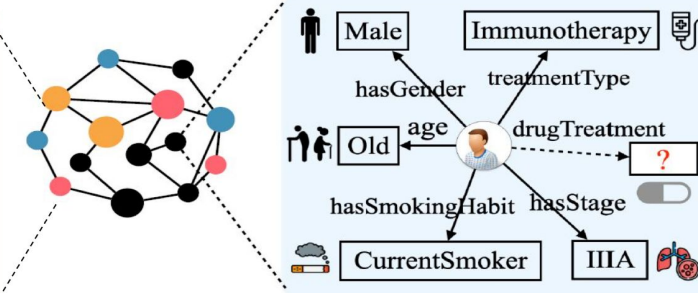
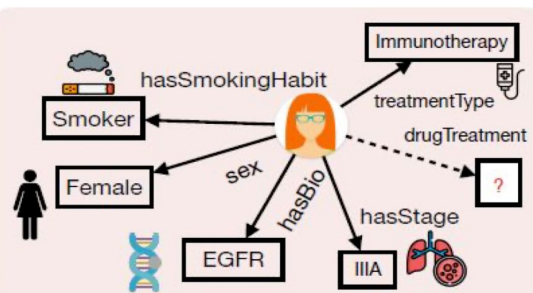
- Address data privacy, sovereignty, and regulations (e.g., GDPR, HIPAA).
- Require transparent frameworks for decision-making and accountability.

TrustKG-Hybrid AI framework to bridge symbolic reasoning and inductive learning and deliver interpretable and user-centric recommendations



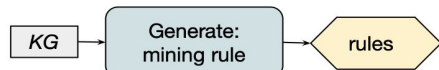
- Semantic Data Integration: Combines biomedical data using modeling, ontology alignment, and semantic reconciliation.
- Knowledge Graph Analytics: Enables link prediction and causal analysis for transparency and interpretability.
- User-Centric Tools: Tools like ORKG Ask provide actionable insights and foster trust.

WISE-Hybrid AI for Accurate and Interpretable Link Prediction in KGs



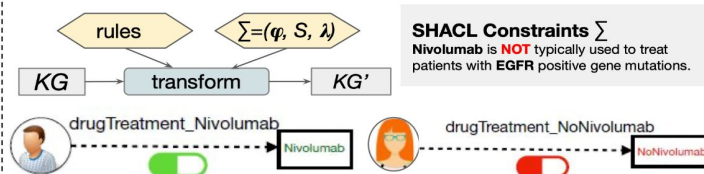
Approaches	Model	Results for KG_3				
		Hits@1	Hits@3	Hits@5	Hits@10	MRR
Baseline 1	TransE	0.000	0.560	0.795	0.943	0.324
	TransD	0.002	0.551	0.690	0.872	0.310
	TransH	0.622	0.864	0.943	0.983	0.756
	RotatE	0.696	0.933	0.969	0.987	0.820
Baseline 2	TransE	0.000	0.713	0.840	0.931	0.376
	TransD	0.008	0.694	0.824	0.935	0.379
	TransH	0.882	0.969	0.997	1.000	0.929
	RotatE	0.864	0.987	0.995	1.000	0.924
Baseline 3	TransE	0.000	0.519	0.747	0.923	0.310
	TransD	0.011	0.551	0.716	0.884	0.322
	TransH	0.596	0.876	0.925	0.977	0.740
	RotatE	0.714	0.941	0.969	0.990	0.829
Baseline 4	TransE	0.000	0.536	0.735	0.931	0.311
	TransD	0.002	0.551	0.733	0.870	0.318
	TransH	0.542	0.849	0.908	0.974	0.702
	RotatE	0.700	0.945	0.972	0.992	0.818
WISE	TransE	0.000	0.760	0.878	0.948	0.388
	TransD	0.013	0.684	0.762	0.884	0.368
	TransH	0.868	0.980	0.994	1.000	0.924
	RotatE	0.887	0.986	0.996	0.998	0.936

a) Symbolic Learning

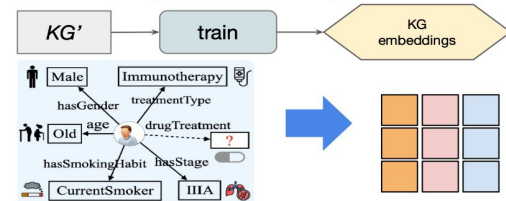


ic:patientDrug(X, Nivolumab):-
 ic:hasStage(X, IIIA),
 ic:treatmentType(X, Immunotherapy).

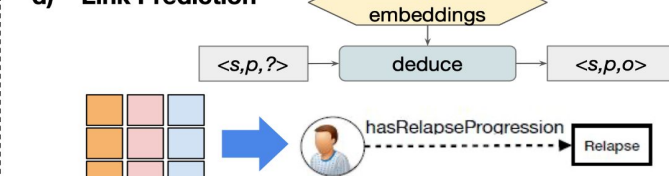
b) KG Validation and Transformation



c) Learning KG Embeddings



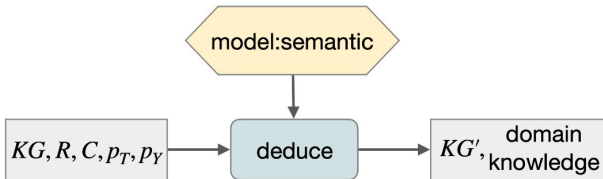
d) Link Prediction



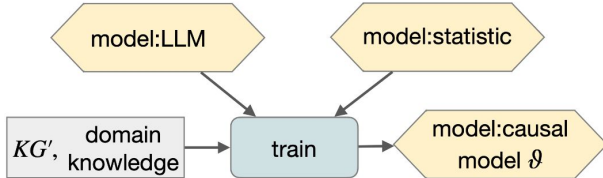
Existing KG embedding models: -)
 impacted by the representation
 of factual statements in KGs
 -) enhanced by explicitly
 expressing valid and invalid links
 according domain-specific
 integrity constraints.

HealthCareAI-Hybrid AI for Accurate Counterfactual Predictions

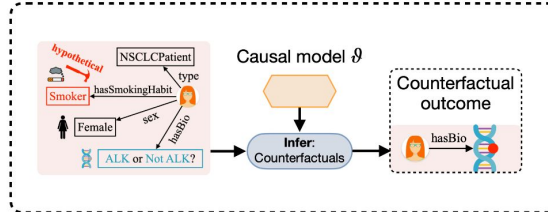
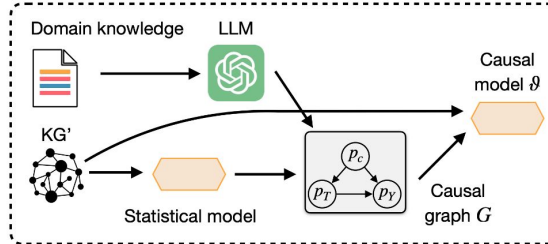
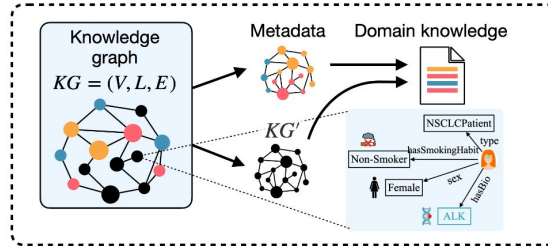
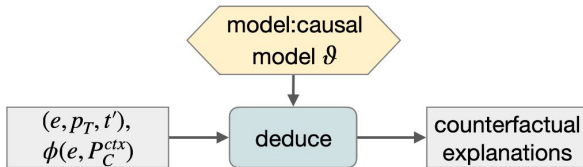
a) Symbolic Reasoning



b) Causal Model Learning



c) Counterfactual Prediction



Expert Causal Graph

Approaches	Model	Results (%) for KG_{real}			
		Jaccard Index	Precision	Recall	F1-Score
Baseline 1	PC	13.3	40.0	16.7	23.5
	FCI	16.0	23.5	33.3	27.6
	GES	6.7	25.0	8.3	12.5
Baseline 2	GPT4 wo domain knowledge	0.0	0.0	0.0	0.0
Baseline 3	GPT4 w domain knowledge	23.5	44.4	33.3	38.1
HealthCareAI	Baseline 3 + PC	30.0	42.9	50.0	46.2
	Baseline 3 + FCI	29.6	34.8	66.7	45.7
	Baseline 3 + GES	25.0	38.5	41.7	40.0

Counterfactual Prediction

Approaches	Model	Results (%) for KG_{10k}			
		Jaccard Index	Precision	Recall	F1-Score
Baseline 1	PC	75.0	100.0	75.0	85.7
	FCI	75.0	100.0	75.0	85.7
	GES	62.5	100.0	62.5	76.9
Baseline 2	GPT4 wo domain knowledge	40.0	66.7	50.0	57.1
Baseline 3	GPT4 w domain knowledge	62.5	100.0	62.5	76.9
HealthCareAI	Baseline 2 + PC	70.0	77.8	87.5	82.4
	Baseline 3 + PC	100.0	100.0	100.0	100.0

-) Compared against statistical methods (PC, FIC, GES) and LLM-based approaches with/without domain knowledge.
-) Achieved higher accuracy by integrating semantics and domain-specific knowledge.
-) Semantic enrichment improved interpretability and aligned predictions with human reasoning.

Open Research Knowledge Graph (ORKG)- represents, organizes, and shares scholarly knowledge in a structured and machine-readable format



Prof. Dr. Sören Auer

Purpose and Vision:

Aims to transform scientific knowledge into a connected, interactive, and navigable structure that enhances understanding and accessibility.

Key Features:

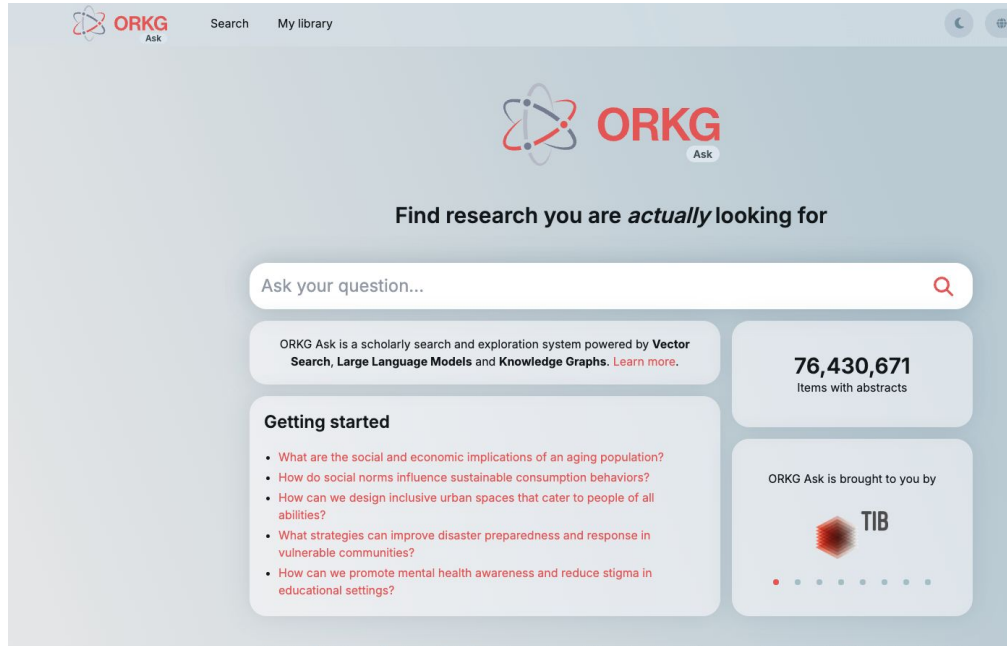
Structured Data Representation: Utilizes a knowledge graph structure to capture and link research outputs for enhanced discoverability.

Semantic Enrichment: Integrates related concepts and research findings to provide contextual and comprehensive insights.

Interactive Interface: Allows researchers to contribute and explore research data through an intuitive user experience.

The screenshot shows the ORKG website interface. At the top, there is a navigation bar with the ORKG logo, a search bar, and options to 'Add new' and 'Sign in'. Below the navigation bar, the main heading reads 'Scholarly Knowledge. Comparable.' followed by a brief description of the ORKG project. A 'Browse by research field' section displays five categories: Arts and Humanities (53 papers, 30 comparisons), Engineering (4052 papers, 330 comparisons), Life Sciences (4923 papers, 237 comparisons), Physical Sciences & Mathematics (16427 papers, 237 comparisons), and Social and Behavioral Sciences (1123 papers, 166 comparisons). The main content area is divided into 'Comparisons', 'Papers', 'Visualizations', 'Reviews', and 'Lists'. The 'Comparisons' section is active, showing a list of research comparisons with details such as title, number of contributions, visualizations, and dates. Examples include 'Dataset in wind energy assessment in Europe', 'Soil Classification and Analysis based on Nutritional Composition', 'Deep Learning Methods for Fake News Detection', 'Branch and Bound algorithms application in works of Librarian scientists', 'Comparative Analysis of ALD-Deposited Films Across Varied Processes', 'Methods for antisocial behavior identification on social networks', and 'Comparison of LLM Hallucination Benchmarks'. A 'Top contributors' section is also visible, listing users like 'Kema score' and 'Vasily Deibert'. At the bottom left, the URL 'https://orkg.org/' is displayed in a large, blue, stylized font.

ORKG Ask- Human-Centric Communication



A scientific search and exploration system providing answers from a database of 80 million full-text scientific publications.

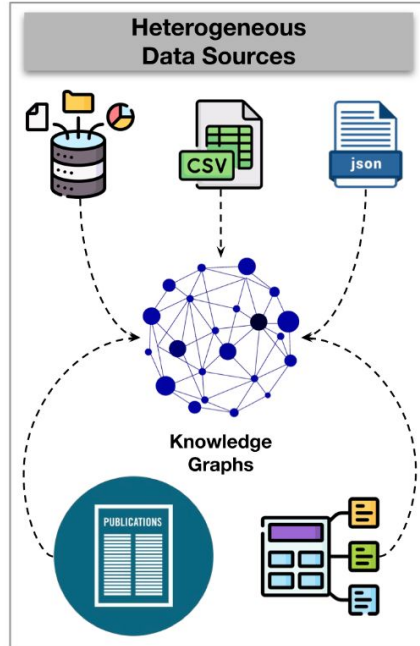
Combines the power of large language models with the structured data of ORKG to deliver fast, informed answers to complex research questions.

<https://ask.orkg.org/>

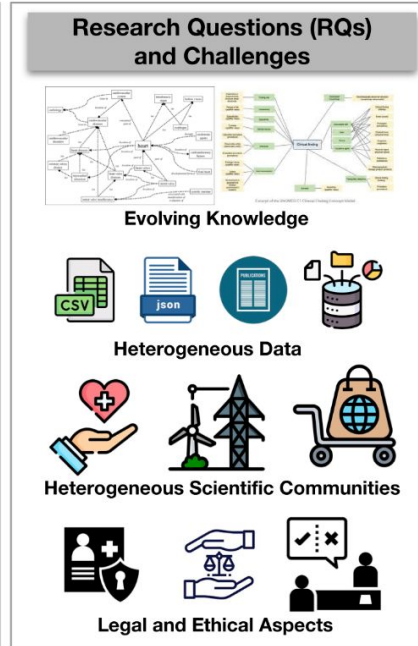
[Example1](#)

[Example2](#)

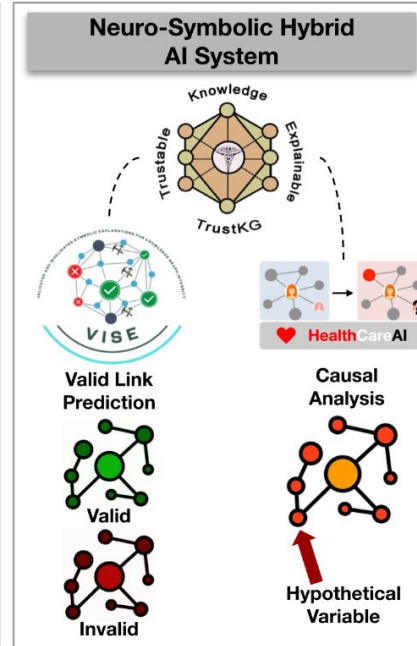
TrustKG- Knowledge Graph-driven AI system to enhance interpretability, transparency, and usability



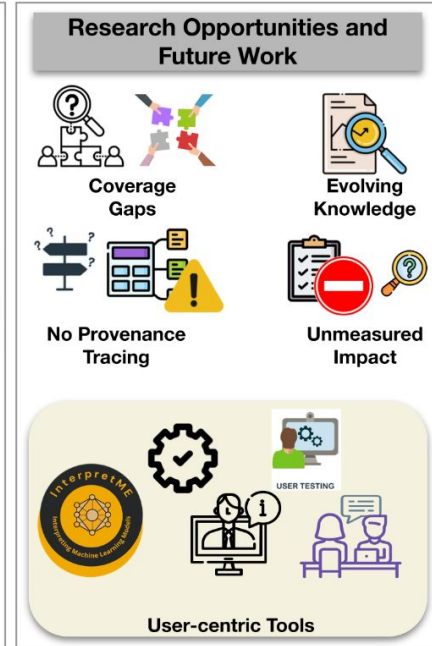
Knowledge Graphs (KGs) present a promising avenue for achieving integration of heterogeneous data sources, thereby facilitating **interoperability** and data exploration.



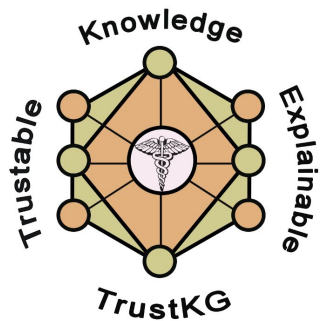
The **research challenges** identified here are critical barriers to integrating KGs with hybrid AI systems.



TrustKG shows how hybrid AI systems like **VISE** and **HealthCareAI** can use KGs in different areas.



By overcoming these limitations and advancing their research, KGs can reach their full potential and use **advanced Hybrid AI** in real-world applications.



Leibniz Programme for
Women Professors



Disha Purohit



Yashrajsinh Chudasama



Hao Huang



The background features several abstract, translucent blue and green shapes that resemble molecular structures or cells. These shapes are scattered across the frame, with some containing small, colorful dots in red, green, and blue. The overall aesthetic is scientific and modern.

Yongxin (Kiki) Kong, *Indiana*
University & Chinese Academy of Sciences
(HRAIit)

HRALit: Publication, funding, and experimental data in support of Human Reference Atlas construction and usage

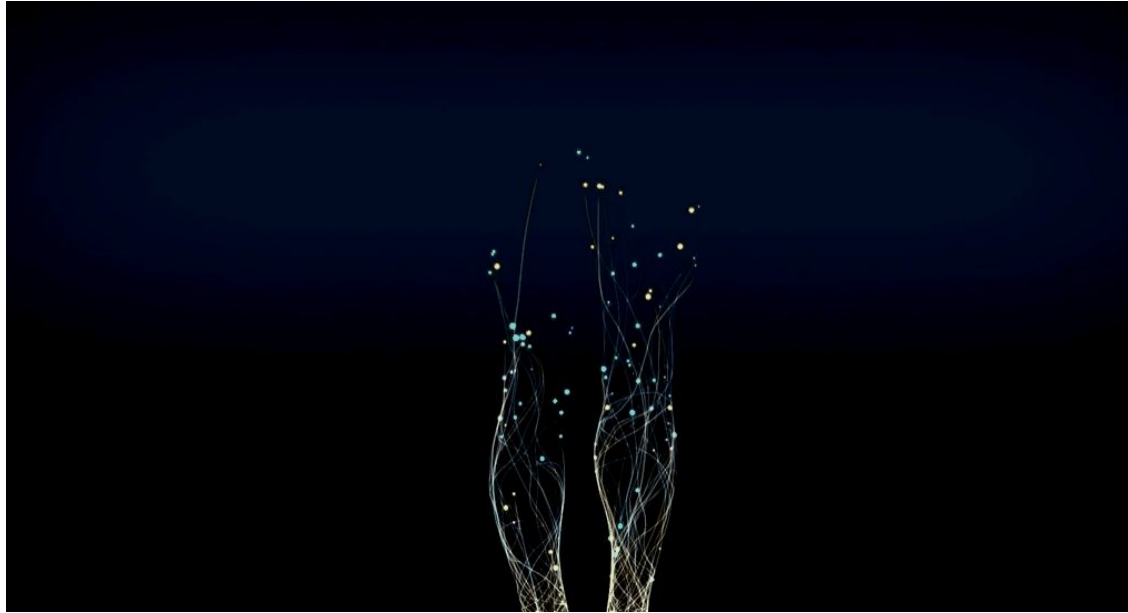
Yongxin (Kiki) Kong

Postdoc, Chinese Academy of Science

Visiting Ph.D., Indiana University



The Human Reference Atlas (HRA) effort aims to map the human body at single cell resolution.

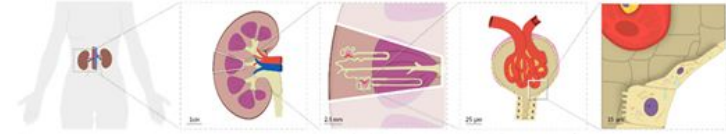
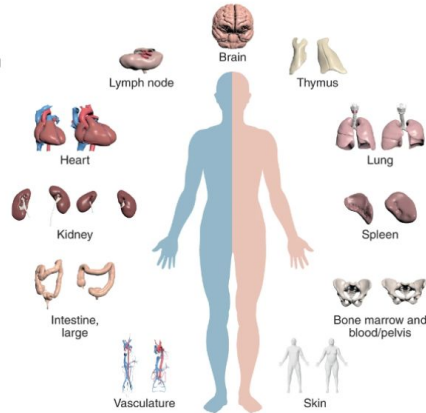
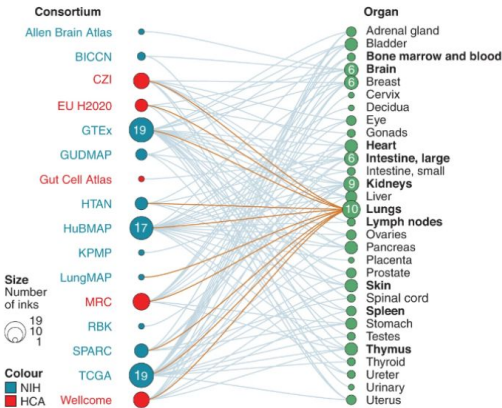


<https://humanatlas.io>



<https://3d.nih.gov/collections/hra>

HRA grows in the number of organs and data types



- | Body | Organ | Functional Tissue Unit | FTU Sub-structure(s) | Cellular |
|---|--|--|--|--|
| <ul style="list-style-type: none"> • Body • Kidney (Left, Right) • Aorta • Renal artery • Renal vein • Ureter | <ul style="list-style-type: none"> • Renal capsule • Renal pyramid • Renal cortex • Renal medulla • Renal calyces • Renal pelvis | <ul style="list-style-type: none"> • Nephron • Renal corpuscle • Proximal convoluted tubule • Loop of Henle • Distal convoluted tubule • Collecting duct | <ul style="list-style-type: none"> • Bowman's capsule • Glomerulus • Efferent arteriole • Afferent arteriole | <ul style="list-style-type: none"> • Parietal epithelial cells • Endothelial cells • Mesangial cells • Podocytes |



17
consortia



250+
experts



1,000+
publications



65
organs



4,659
anatomical
structures



1,311
cell types



2,024
biomarkers



21
organ mapping
antibody
panels

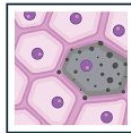


22
functional
tissue units

Many high-quality experimental datasets are becoming available



Human BioMolecular Atlas Program



Cellular Senescence Network



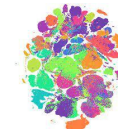
Genotype-Tissue Expression



Kidney Precision Medicine Project



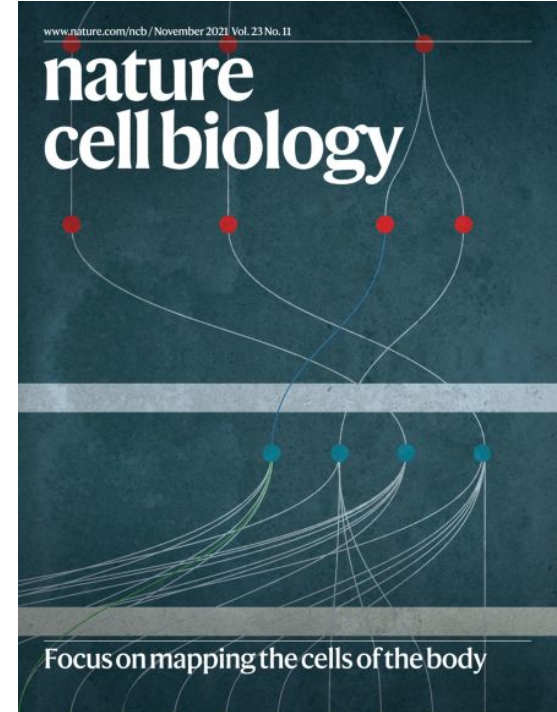
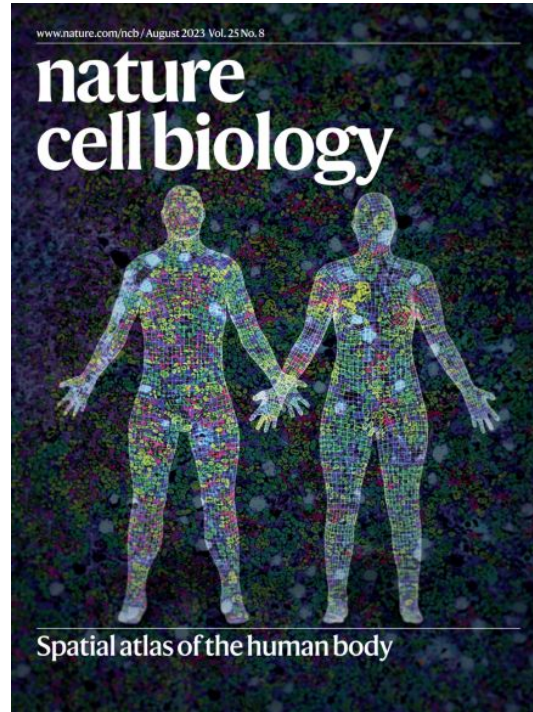
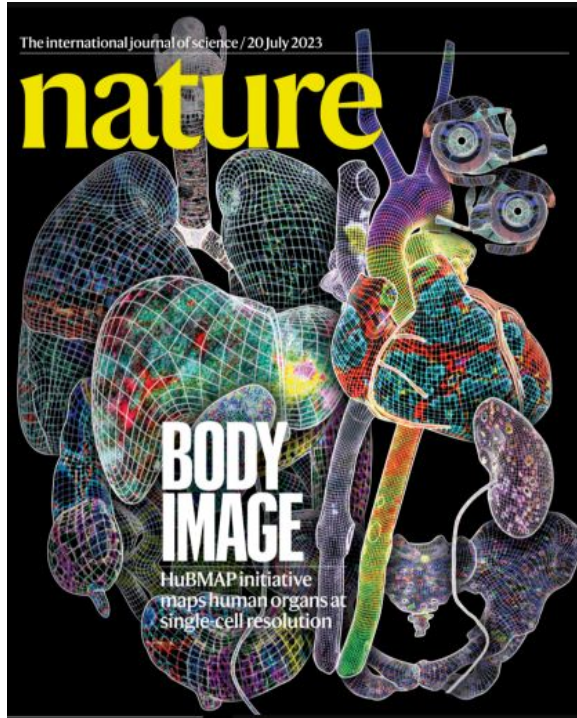
GenitoUrinary Developmental Molecular Anatomy Project



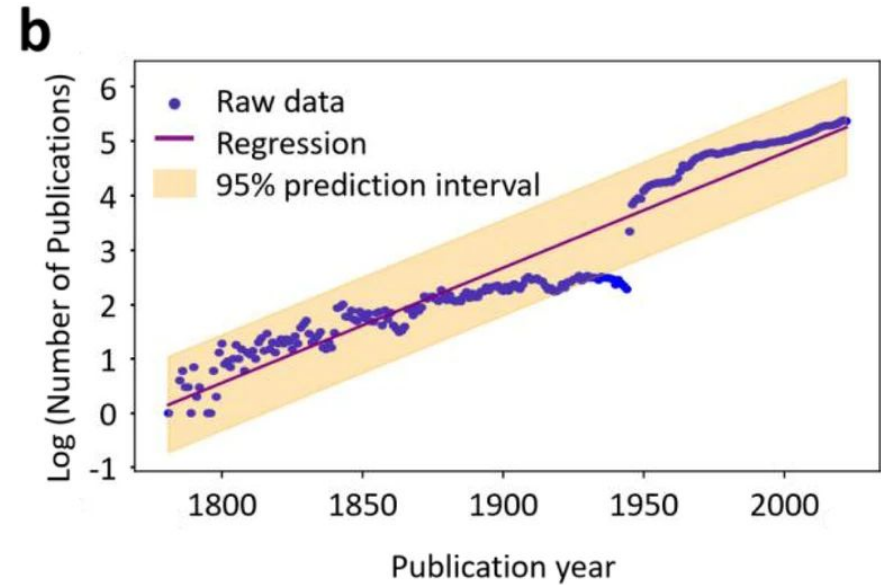
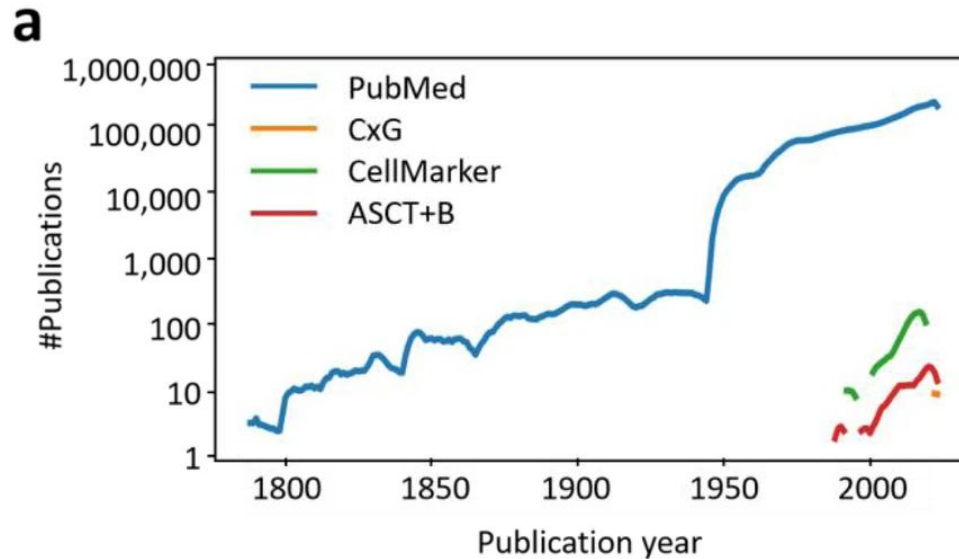
CZ CELLxGENE

Searching for data across portals is difficult

HRA relevant data is published in scholarly papers



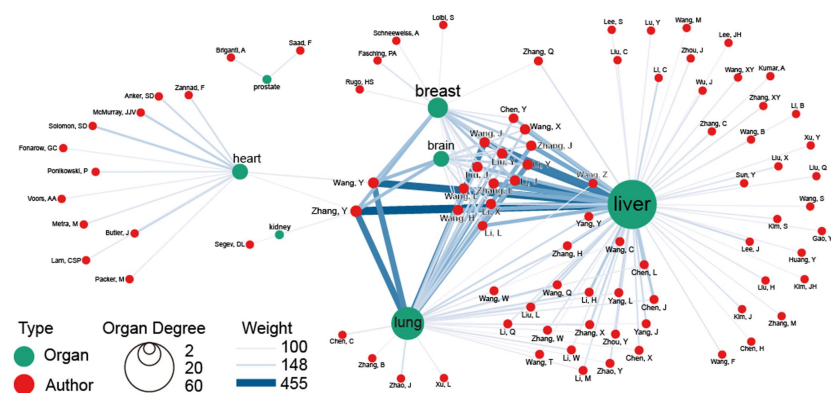
The number of publications increases exponentially



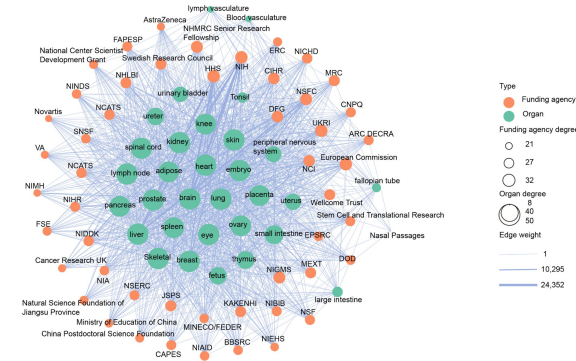
Scholarly publications evidence for HRA



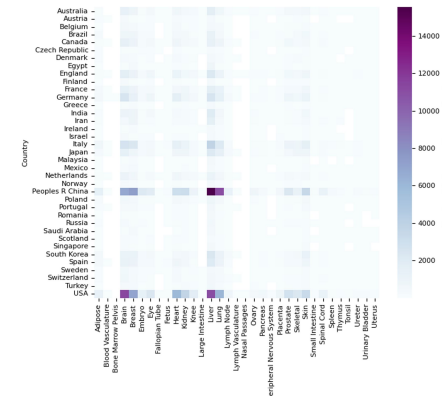
Geospatial layout of the coauthor network



Bimodal network of highly cited authors and the organs they study.



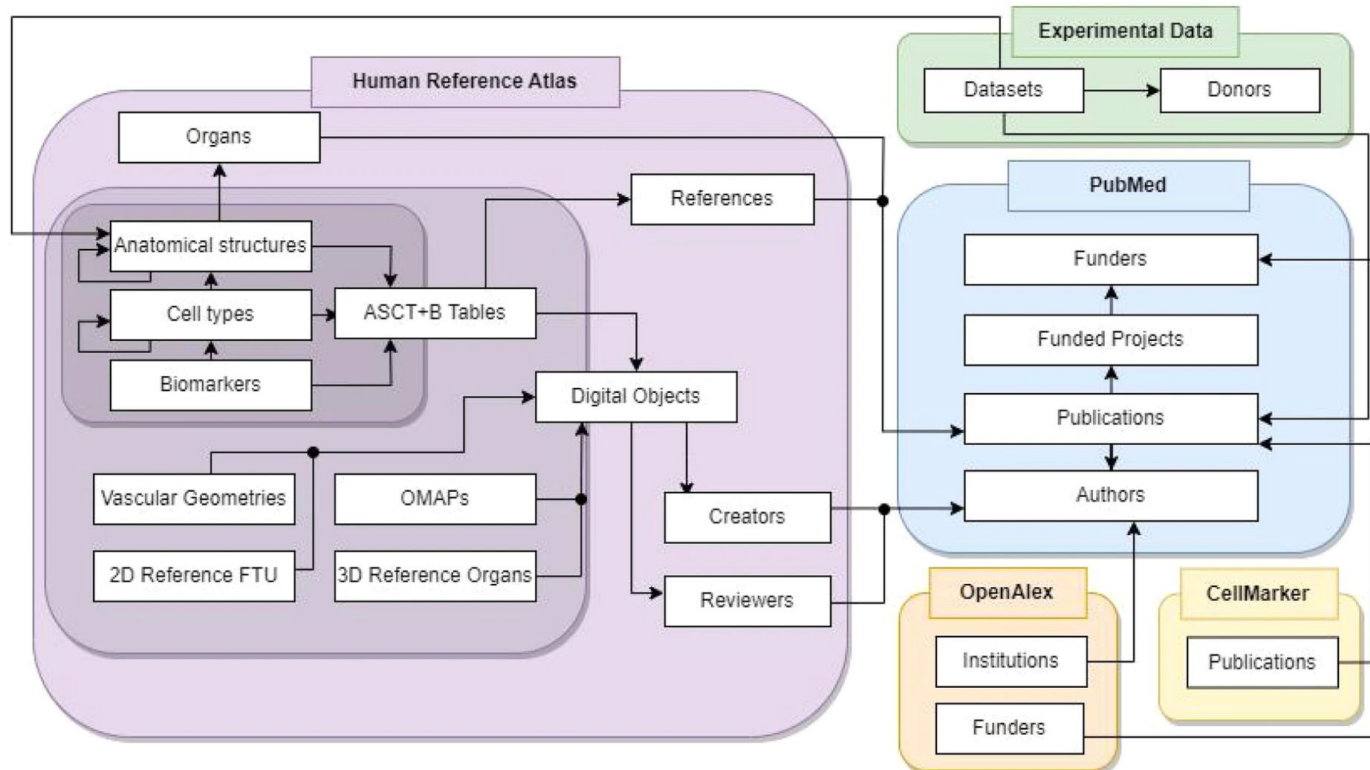
Bimodal network of 32 organs and top 50 funding agencies most often listed



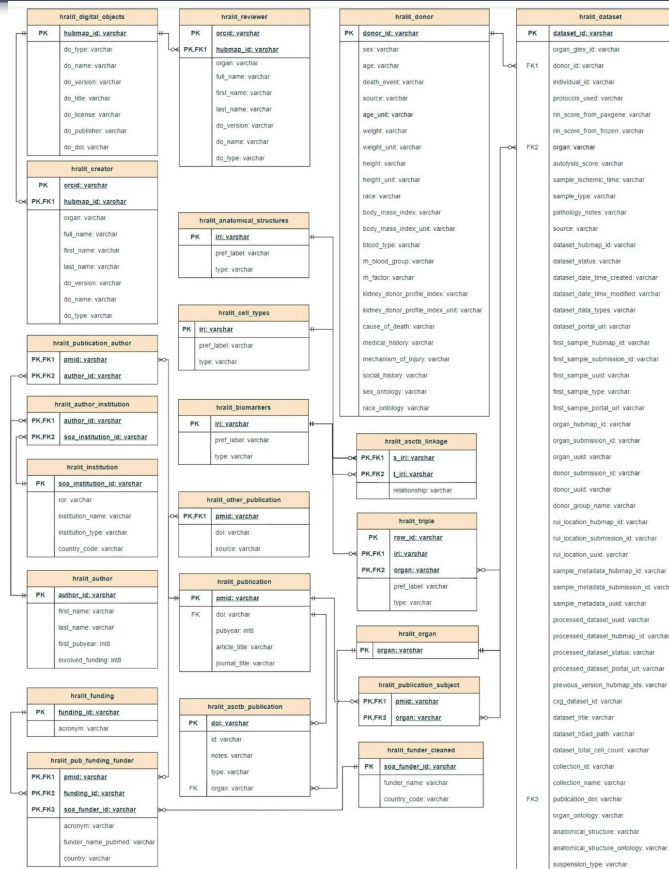
Heatmap of authors per country per organ

Kong Y, Daiya V A, Börner K. *Quantitative Science Studies*, 2024

Overview of the HRAlit database



Entity relationship diagram of the HRAIt database



The HRAlit database SQL file and all tables in CSV format are at Figshare



Browse

Search on figshare...



Log in Sign up





DATASET

hralit_anatomi... (368.51 kB)  





DATASET

hralit_asctb_link... (2.41 MB)  



DATASET

hralit_asct_pu... (178.92 kB)  





DATASET

hralit_biomark... (113.68 kB)  



.GZ

hralit_author.csv... (3.92 MB)  

1/3





DATASET

hralit_cell_types... (87.95 kB)  



.GZ

hralit_author_ins... (3.23 MB)  



DATASET

hralit_creator.csv (55.32 kB)  





DATASET

hralit_digital_ob... (75.52 kB)  



DATASET

hralit_dataset.csv (3.49 MB)  

Give Feedback

Human Reference Atlas Literature (HRAlit) Database

Cite

Download all (1.14 GB)

Share

Embed

+ Collect

Version 2  Dataset posted on 2024-02-01, 08:44 authored by [Yongxin Kong](#), Katy Börner

USAGE METRICS 

464
views

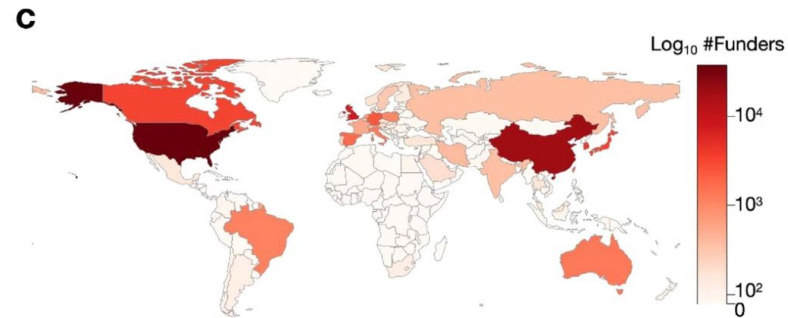
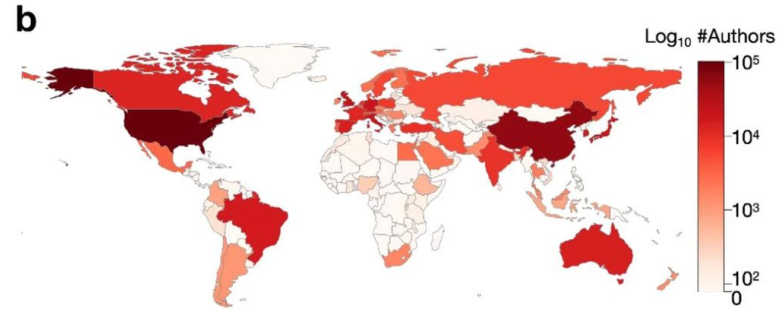
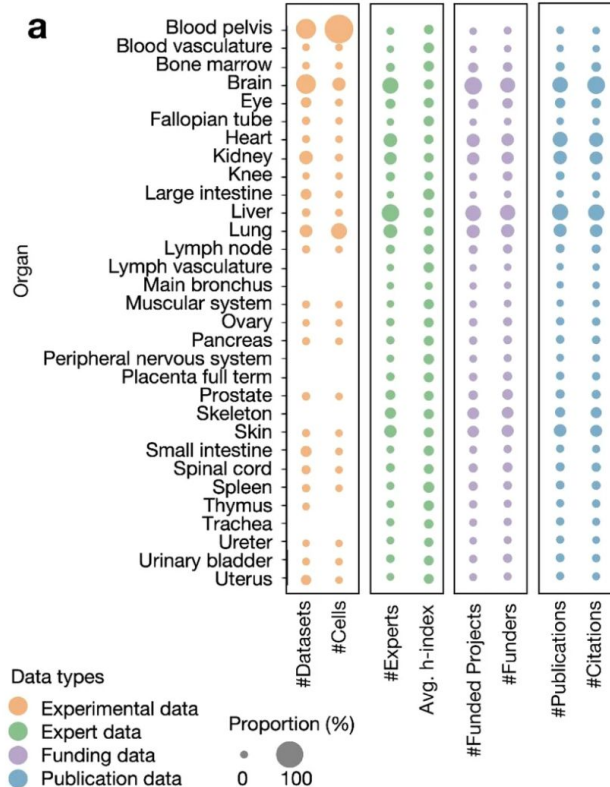
334
downloads

1
citations 

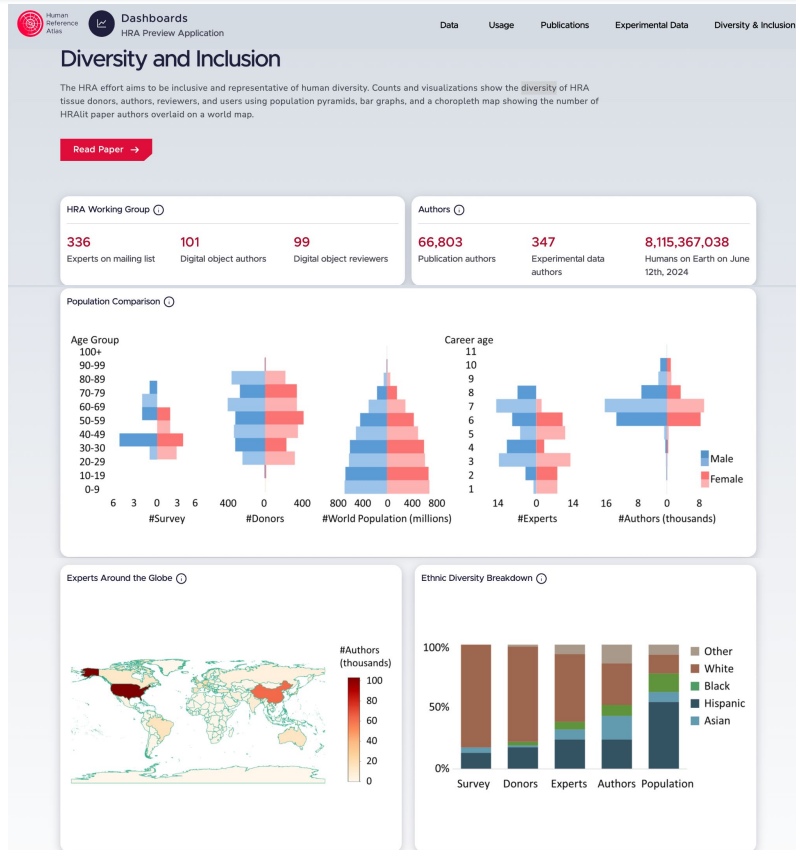
The Human Reference Atlas literature (HRAlit) database, with 23 tables with 21,703,812 records including 7 junction tables with 13,042,188 relationships and a total size of 1.56 GB, is available in SQL format together with tables in CSV format.

Give Feedback

Distribution and choropleth map from HRAlit



HRA diversity and inclusion



Acknowledgements

Principal Investigator:



Katy Börner

Lab:



Funders:



SenNet



NIDDK



The background of the slide features a complex, abstract simulation of particles. It consists of several large, irregular, semi-transparent blue and light green shapes that resemble clusters or aggregates. These shapes are filled with numerous small, multi-colored dots in shades of red, green, and blue. The overall appearance is that of a dynamic, multi-scale system, possibly representing a biological or physical process like protein folding or material aggregation. The text is overlaid on the left side of this simulation.

Oliver He, *University of Michigan*

Ontology: Foundation of Precision Health Data Standardization & Artificial Intelligence

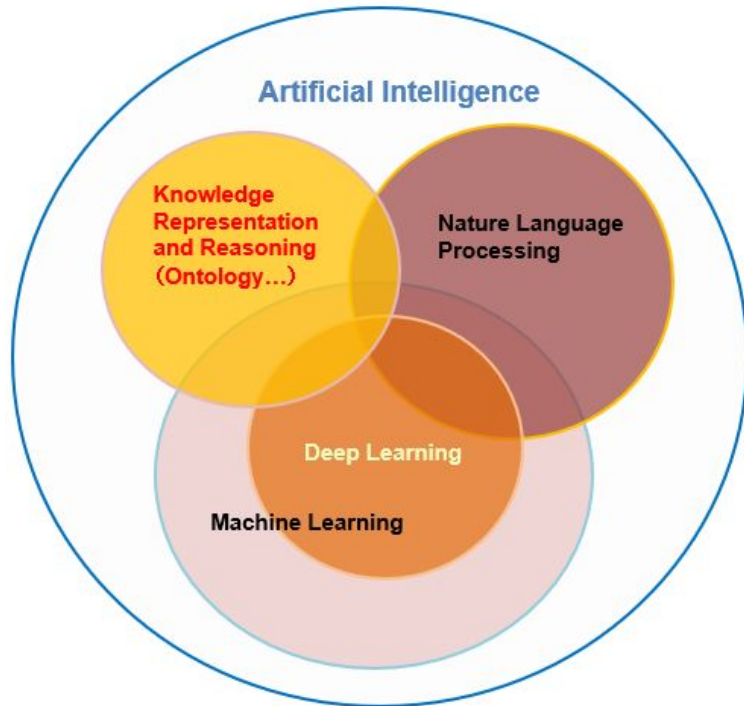
Yongqun “Oliver” He

University of Michigan Medical School
Ann Arbor, MI, USA.



3 Complementary AI Fields: KRR, ML, NLP

(Ontology is a major part of KRR)



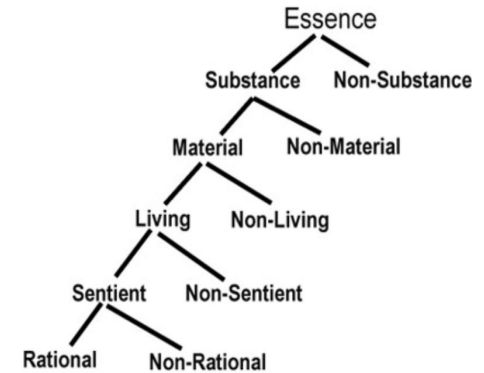
- **KRR**: Knowledge Representation & Reasoning
 - **Ontology** is a major approach in KRR
- **ML**: Machine Learning
- **NLP**: Natural Language Processing
 - Large Language Models (**LLM**): emerging NLP

Ontology: Originated from Philosophy and Taxonomy

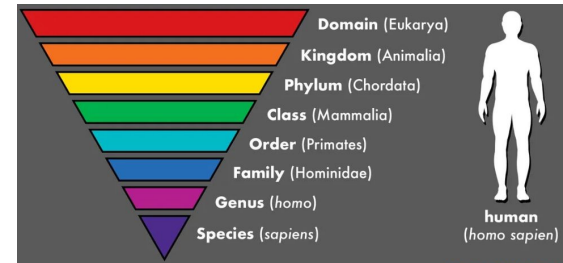
Onto: Being, existence, reality

Definition in philosophy: Ontology is the philosophical study of the nature of being, becoming, existence and/or reality, as well as the basic categories of being and their relations.

- It was called “first philosophy” by Aristotle (384–322 BC) in Book IV of his *Metaphysics*.
- The term “ontology” (or ontologia, “science of being”) was coined in 1613, independently, by two philosophers Rudolf Göckel (Goclenius) in his *Lexicon Philosophicum* and Jacob Lorhard (Lorhardus) in his *Theatrum philosophicum*.
- Stages of “ontology” development as science:
 - Porphyrian tree or Tree of Porphyry (234 – 305 AD)
 - Taxonomy, e.g., Taxonomy of Linnaeus (1707 - 1778)

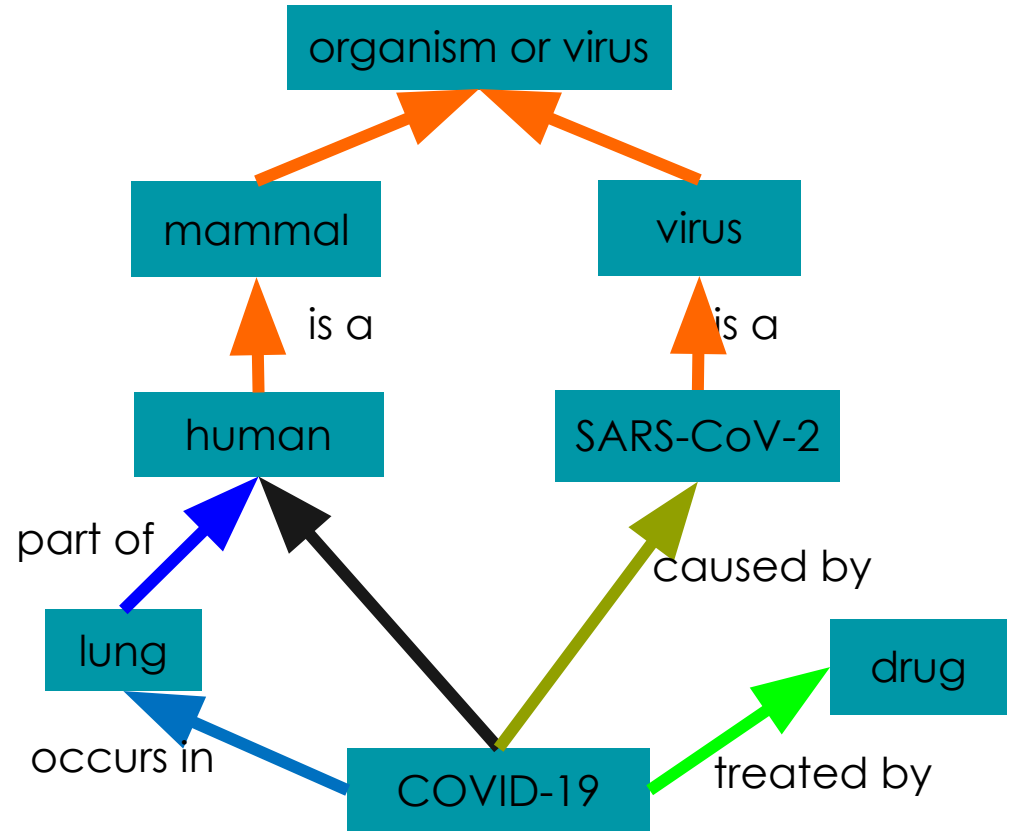


Porphyrian tree



Ontology in IT

- In IT era, ontology is human- and **computer-interpretable** representation of **entities** and the **relations** among entities in a specific domain.
 - A complex, standardized, and integrative network
- Foundation of AI knowledge representation and reasoning
- Support data/knowledge standardization, annotation, integration, and reasoning.



Ontology: Language of AI - Connecting machines & humans

Ontology = Entity terms (controlled vocabulary)
+ Relations (semantics)

Semantic: of, relating to, or arising from the **meanings** of words

“The Semantic Web is an extension of the current web in which information is given well-defined **meaning**, better enabling computers and people to work in cooperation.”

- Tim Berners-Lee, Inventor of WWW

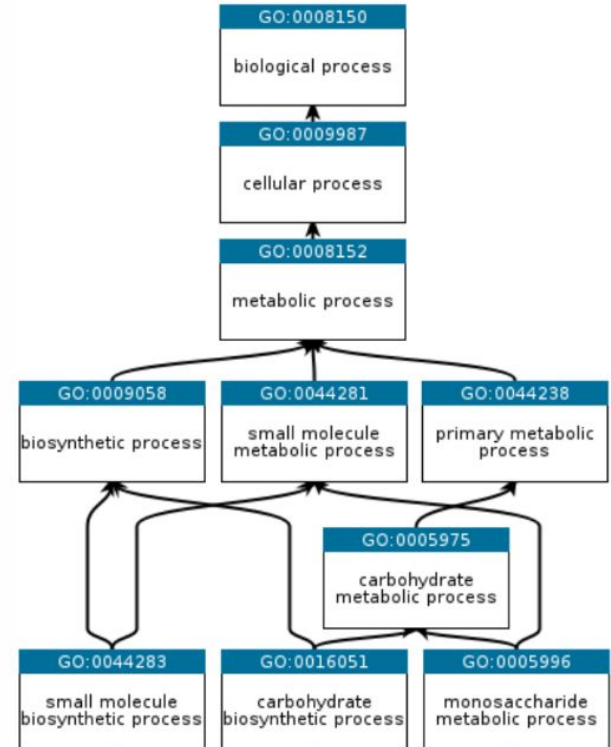
https://en.wikipedia.org/wiki/Semantic_Web



Tower of Babel

Gene Ontology (GO): Critical to gene/genome annotation & gene enrichment analysis, etc.

- Established in 1998 by a consortium studying the genomes of 3 model organisms: *D. melanogaster* (**fruit fly**), *M. musculus* (**mouse**), and *S. cerevisiae* (baker's **yeast**).
- **Goal:** Standardize/unify representation of gene functional annotation across databases and organisms.
 - 3 branches: Biological Process, Cellular Component, Molecular Function
- Critical to gene annotation and analysis.
- Stimulate more ontology development.





Open Biological and Biomedical Ontology (OBO) Foundry

- International open initiative since 2006
- A collection of orthogonal reference ontologies in biological and biomedical domain.
 - Gene Ontology was the first one joining.
 - Other examples: Cell Ontology (CL), Mammalian Phenotype Ontology (MP), Human Phenotype Ontology (HPO), **Vaccine Ontology (VO)**, etc.
- Each is committed to a set of **principles** for the best practices in ontology development.

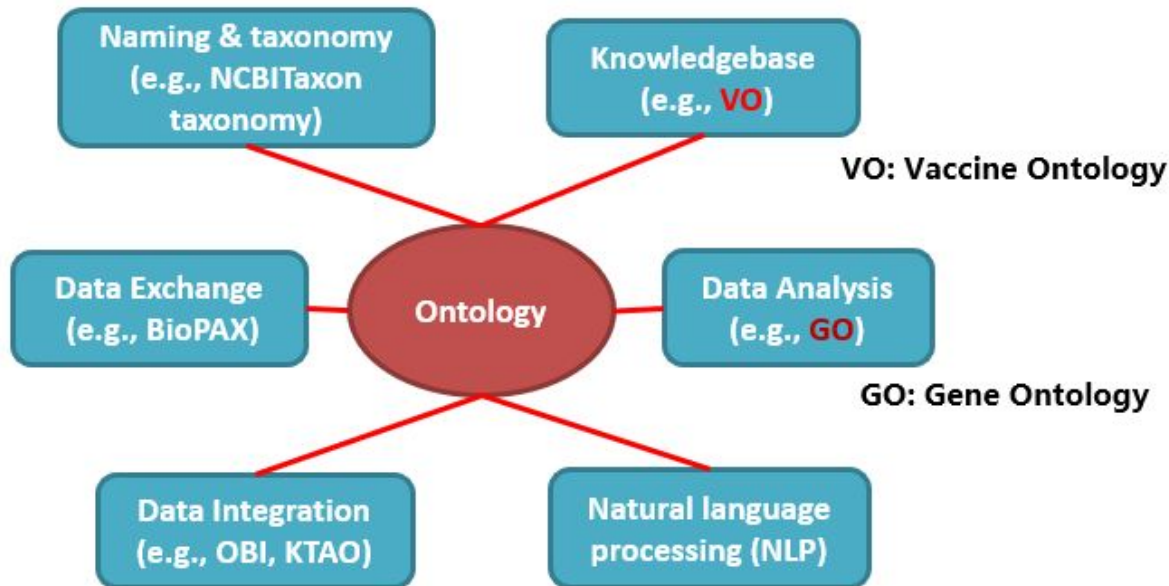


OBO Principles:

- Open (principle 1)
- Common Format (principle 2)
- URI/Identifier Space (principle 3)
- Versioning (principle 4)
- Scope (principle 5)
- Textual Definitions (principle 6)
- Relations (principle 7)
- Documentation (principle 8)
- Documented Plurality of Users (principle 9)
- Commitment To Collaboration (principle 10)
- Locus of Authority (principle 11)
- Naming Conventions (principle 12)
- Notification of Changes (principle 13)
- Maintenance (principle 16)
- Responsiveness (principle 20)

Applications of Ontology and Semantic Web

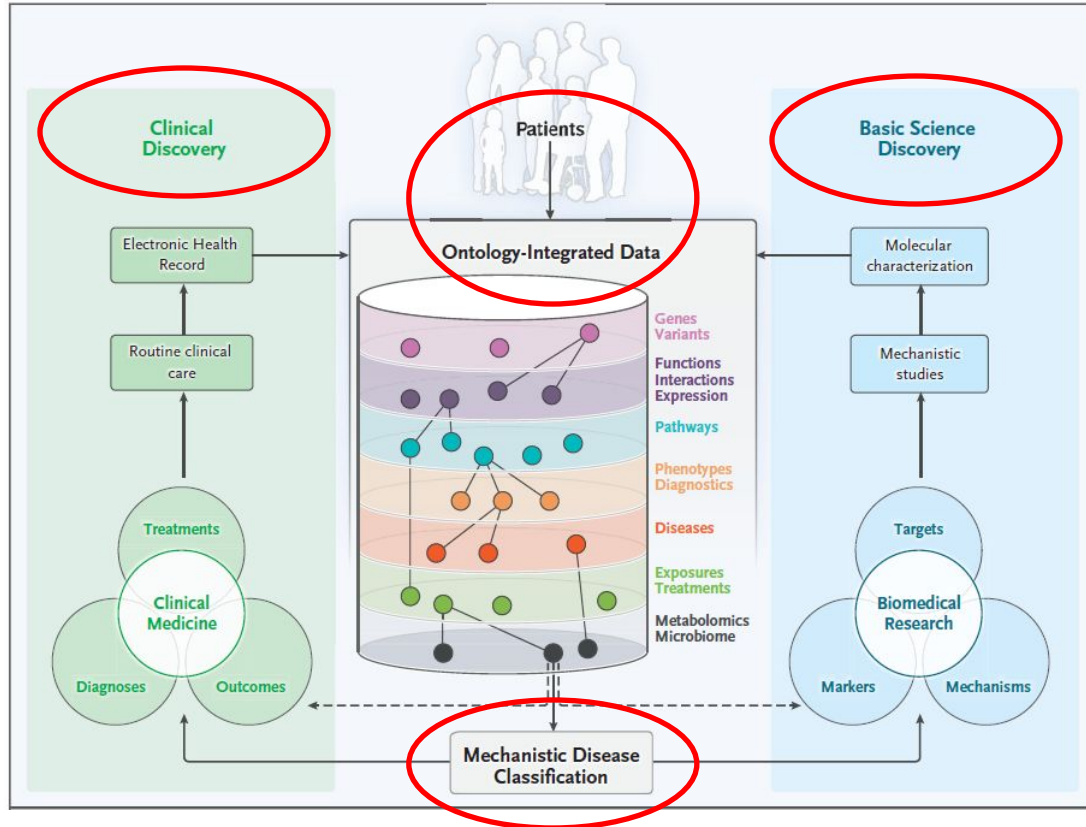
Ontology = Terms (controlled vocabulary)+ Relations (semantics)



OBI: Ontology for Biomedical Investigations

KTAO: Kidney Tissue Atlas Ontology

Ontology: Foundation of Precision Medicine

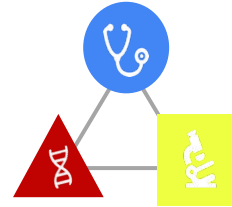


Integrating the two streams of data (**clinical** and **basic** science observations) enables more refined and dynamic classification of disease across many data types

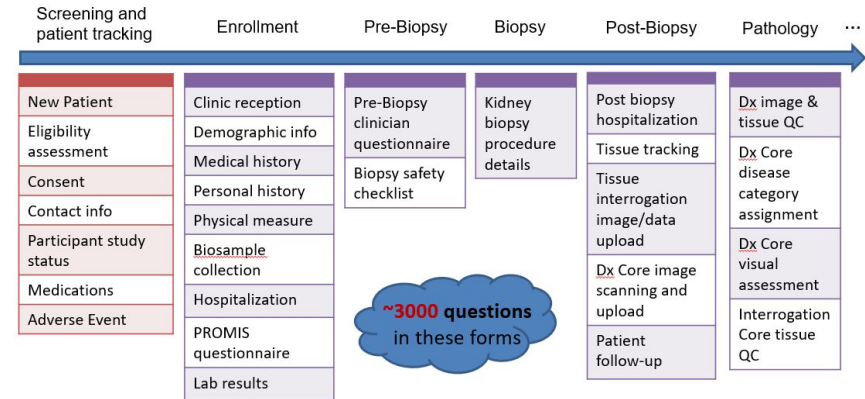
Citation: Haendel MA, Chute CG, Robinson PN. **Classification, Ontology, and Precision Medicine.** *N Engl J Med.* 2018 Oct 11; 379(15): 1452-1462.

KPMP: Opportunities and Challenges

- Initiated 2017, Kidney Precision Medicine Project (**KPMP**), funded by NIH-NIDDK
 - Only human studies, no lab animals.
- Over 20 universities / institutes
- **Goals:**
 - Build a **kidney tissue atlas** that links *clinical* phenotypes, cells, *molecules*, pathways, and *pathology* together.
 - Understand and treat **human kidney diseases** – Acute Kidney injury (**AKI**) and Chronic Kidney Disease (**CKD**)
- **“Big data” challenge:** integration & analysis



- **Clinical**
- **Molecular**
- **Pathology**



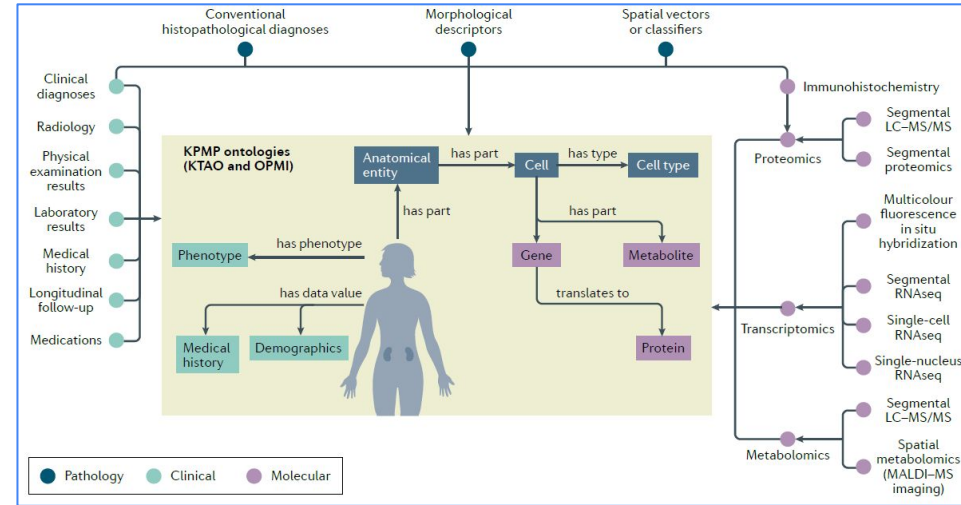
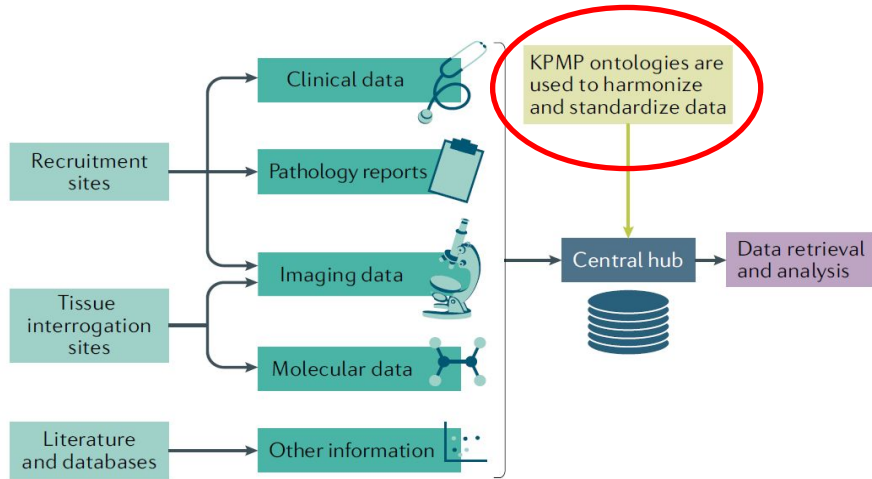
38 KPMP Case Report Forms (CRFs)

Two ontologies for KPMP

- Two community-based KPMP ontologies:
 - **KTAO: Kidney Tissue Atlas Ontology** – It's more about [kidney knowledge](#)
 - **OPMI: Ontology of Precision Medicine and Investigations** – Standardizes [data and metadata](#) types in and beyond KPMP.
 - Kidney-related info in OPMI is imported back to KTAO.
- Interoperable ontology development strategies
 - Follow Open Biomedical Ontology (**OBO**) **principles**: Openness, collaboration, etc.
 - >150 OBO library ontologies: **non-redundant, interoperable**
 - **Reuse/align/integrate** existing ontologies: UBERON anatomical entity, HPO (Human Phenotypes), GO, CL (Cells), OBI (Biomedical Investigations), ...

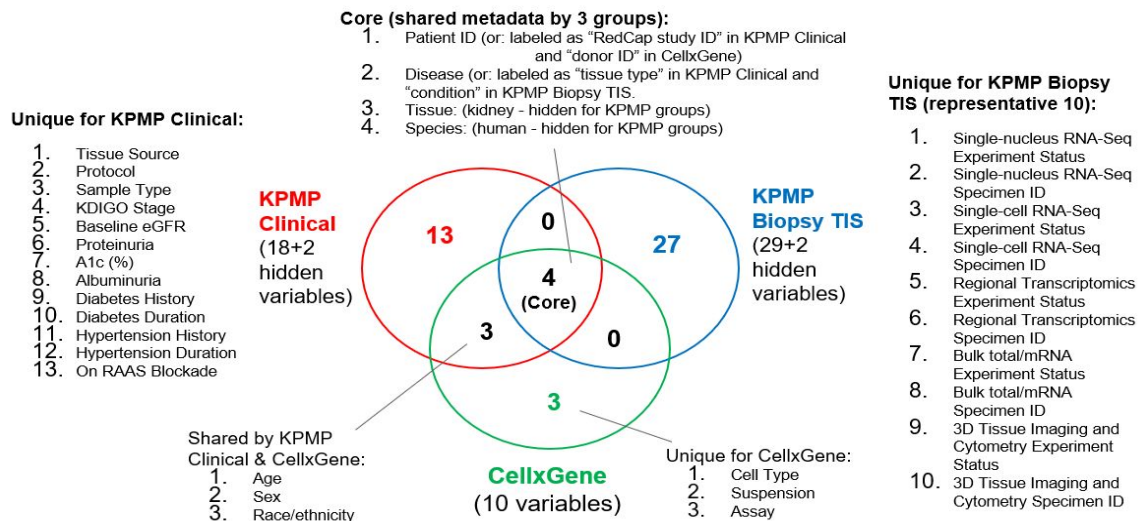
Ref: Ong E, Wang LL, Schaub J, O'Toole JF, Steck B, Rosenberg AZ, Dowd F, Hansen J, Barisoni L, Jain S, de Boer IH, Valerius MT, Waikar SS, Park C, Crawford DC, Alexandrov T, Anderton CR, Stoeckert C, Weng C, Diehl AD, Mungall CJ, Haendel M, Robinson PN, Himmelfarb J, Iyengar R, Kretzler M, Mooney S, and He Y, for the Kidney Precision Medicine Project. Modeling Kidney Disease Using Ontology: Perspectives from the KPMP. *Nature Review Nephrology*. 2020 Nov;16(11):686-696. PMID: 32939051.

Ontology critical to KPMP big data integration and analysis



Ong, et al., *Nature Review Nephrology*, 2020

Integrate KPMP, HuBMAP, and CellxGene data using ontology



- **HuBMAP** (Human BioMolecular Atlas Program) focuses on reference human body.
- **KPMP** focuses on diseased kidney
- Even so, the results are not naturally integrated
- We proposed an interoperable ontology **“Precision Medicine Metadata Ontology (PMMO)”** to harmonize and integrate the data.

He Y, et al. AMIA 2024 full-length paper, oral presentation, Best Paper Award

HuBMAP Hackathon: Integrating KPMP & HuBMAP data by ontology

- **Goal:** Harmonize and integrate KPMP and HuBMAP data to more efficiently address scientific questions
- <https://github.com/hubmapconsortium/hra-hubmap-kpmp-integration>
- **Key process:** metadata harmonization using ontology
 - PMMO: Precision Medicine Metadata Ontology
- **Use case study:**
 - **SPP1:** A biomarker that differentiates healthy from AKI

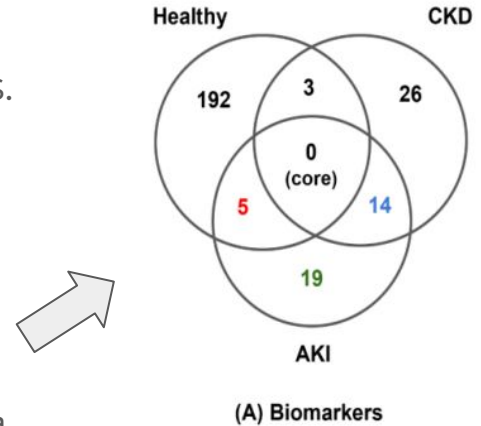


Team members:

- Bruce Herr (IU)
- Yichao Chen (PSU)
- Leo Yeh (UM)
- Ruopeng Wu (UM)
- Oliver He (UM)

Biological Insights:

- **Theme:** Compare **gene biomarkers** between healthy and disease → leverage the info to **cell type** and **anatomy** levels.
 - This aligns with HubMAP **Anatomical Structures, Cell Types, and Biomarkers (ASCT+B) Tables**.
 - HuBMAP has collection of healthy kidney biomarkers
 - KPMP have collections of diseased kidney biomarkers.
- At **gene biomarker** level:
 - 26 AKI biomarkers were found, 5 also shared with healthy kidney
 - Previously, we started with KPMP/cellxgene data.
 - For the Hackathon, we focus on HuBMAP data and merged earlier data
 - **Hypothesis:** An AKI/healthy gene biomarker(s) may have differential gene expression profiles in AKI patients vs healthy human subjects.
- At **cell** level:
 - Many biomarkers are for specific cell types. By analyzing the cell type specific biomarkers, we can indirectly find the cell type expression.
- At **kidney anatomical structure** level:
 - Specific cell types exist in specific regions. Through the **chain of biomarker-cell-Anatomy**, we can infer specific kidney region activities through the gene biomarker expression.
 - HuBMAP/KPMP histological image data can also be used later.

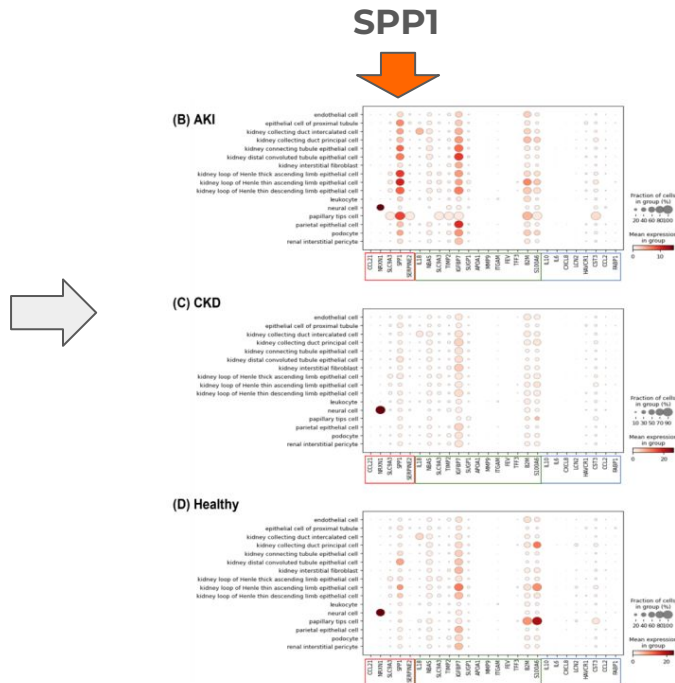


<https://www.biorxiv.org/content/10.1101/2024.04.01.587658v1.full.pdf>

Note: Paper also presented in **AMIA 2024 Annual Symposium**, → **Best Paper Award**.

SPP1: A biomarker that differentiates healthy from AKI

- **SPP1**: Secreted Phosphoprotein 1
 - <https://www.genecards.org/cgi-bin/carddisp.pl?gene=SPP1>
 - Key gene in lymph node metastasis and cancer, but its role in kidney still relatively unclear.
- **Earlier** we used KPMP and CellxGene data:
 - SPP1 significantly differed in gene expression in AKI and healthy groups
- **Now**:
 - **HuBMAP data** is added and merged with KPMP/cellxgene data.
 - Results so far:
 - Extracted SPP1 from KPMP and HuBMAP
 - Differential gene expression profiles found.
- More work ongoing. Ontology level integration as well

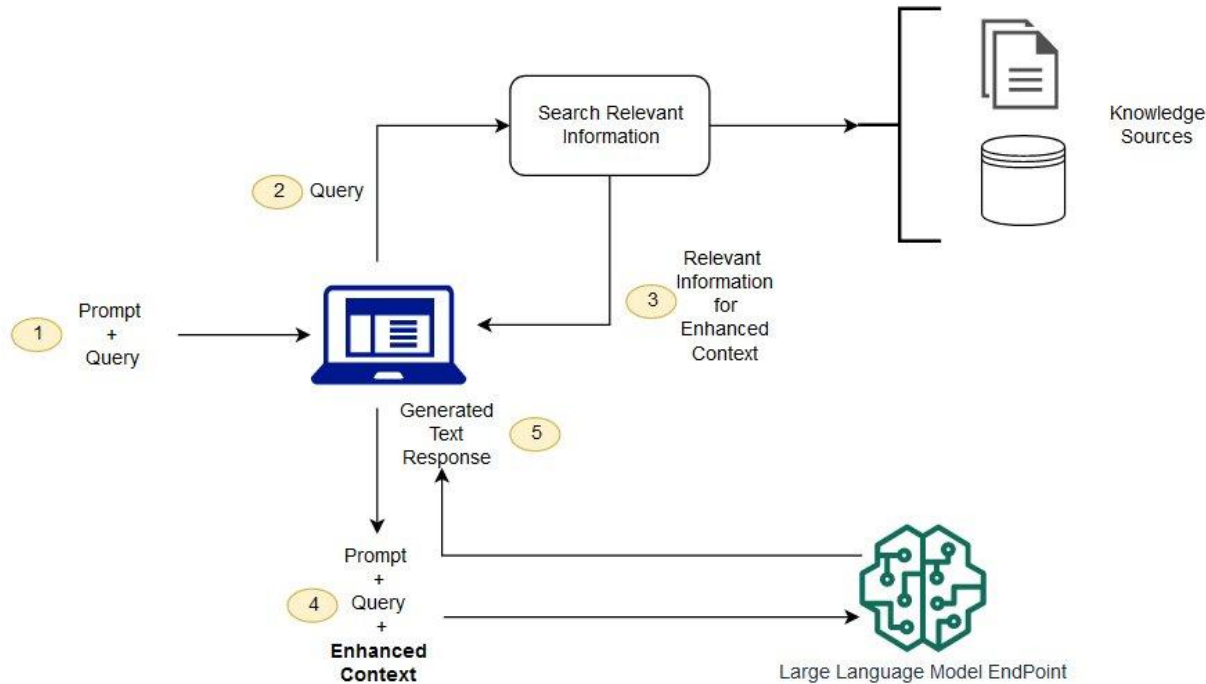


<https://www.biorxiv.org/content/10.1101/2024.04.01.587658v1.full.pdf>

Ontology-Based Knowledge Graph

- **Knowledge graph examples :**
 - Google Knowledge Graph
 - COVID19-KB
- **Two methods of generating knowledge graph:**
 - Triplestore knowledge graph, using tools such as Virtuoso
 - Property knowledge graph, using tools such as Neo4j and GraphDB
- **Ontology role on knowledge graph (KG):**
 - Standardize basic KG framework, including data types & Semantic Relations
 - Computer-understandable knowledge, which can be directly used.
 - Ontology can be used to annotate data

Ontology-Based Knowledge Graph as RAG to Enhance LLM



RAG: Retrieval Augmented Generation

- Optimize LLM output by referencing an authoritative knowledgebase (KB) outside of its training data.
- Such **RAG KB can be generated with ontology support.**

<https://aws.amazon.com/what-is/retrieval-augmented-generation/>

Example: Ontology-supported knowledge RAG for LLM extraction of itemized vaccine information (e.g., vaccine names, types, vaccine antigens, host responses, & experimental factors)

VIOLIN manually collected ~5,000 papers and web links for ~4,700 vaccines and 1600 vaccine antigens as of Nov 1, 2024, since 2007.

- So **~300 per year manually** annotated papers/links.
- Vaccines antigens used as gold standard for vaccine design.

Ontology RAG LLM would help

- Preliminary study approved it.
- **GOAL: 6,000 per paper per year (20-fold more productive).**



Select LLM: llama3-70b-8192 *Llama3 70b used*

Select Embeddings: nomic-embed-text

Add URL to Knowledge Base: *Embedding method*

Add a PDF: *KB website can be added here*

fbioe-11-1121074.pdf 3.8MB *Article PDF loaded here*

Upload JSON for Knowledge Graph

vo_json 9.1MB *VO added as RAG*

Standardized Annotations: *Itemized vaccine info extracted from PDF by RAG LLM*

Vaccine Type:
Conjugate Vaccine (VO:0000163, SubClass Of: "vaccine type", LABEL: "conjugate vaccine")

Vaccine Formulation:
Bioconjugate of Cholera Toxin B Subunit (CTB) with O-Polysaccharides (OPS) from Brucella abortus (VO:0005431, Annotation Property: "vaccine antigen annotation", definition: "An annotation property that represents the annotated information about the antigens used for a specific vaccine.")

Host Species Used as Laboratory Animal Model:
Mus musculus (BALB/c mice) (VO:0000281, LABEL: "laboratory animal model", definition: "A model organism used in scientific research, often used to study human diseases or to test new treatments.")

Experiment Used to Investigate the Vaccine:
Animal Experiment (VO:0000497, LABEL: "animal experiment", definition: "An experiment that uses animals as the subjects, often to study the effects of a vaccine or drug.")

Additional Annotations:

- **Immunization Protocol:** Immunization followed by challenge with lethal and non-lethal doses of Brucella abortus A19 strain (VO:0000503, LABEL: "immunization protocol", definition: "A protocol that outlines the process of immunizing an organism against a specific disease or pathogen.")
- **Immune Response Assays:** Specific antibody production and protective efficacy were assessed, including bacterial load reduction in the spleen and survival rate after lethal

Run ID: d5211256-8d4c-4baa-b6ba-c3c...

Your message

Use case demo: More collected and annotated **vaccine antigens** would serve as **gold standard** for enhanced vaccine design.

Discussion

- **How can ontologies improve AI, & how can AI improve ontologies?**
 - Ontology provides structured knowledge that enhances AI
 - Ontology-supported KG as RAG → improve AI
 - AI can identify new ontology terms/ relations and improve ontology applications
- **How can retrieval-augmented generation (RAG) specifically help with improving ontologies?**
 - Ontology-supported KG as RAG
 - RAG can retrieve docs/data and generate new ontology terms/relations.
- **How can AI chatbots (text-based) accurately represent ontologies (graph-based)?**
 - Chatbots can query ontology-converted KG or triple store
 - Incorporate precomputed inferences into chatbot database

Acknowledgements

U. of Michigan

- Jie Zheng
- Anthony Huffman
- Asiyah Yu Lin
- Leo Yeh
- Edison Ong
- Laurel Li
- Michael Cooke
- Ruopeng Wu
- Oliver He
- Nikki Bonevich
- Jennifer Schaub
- Matthias Kretzler

Penn State

- Yichao Cheng

KPMP

- Jimmy P. Phuong
- Sean Mooney
- Jonathan Himmelfarb
- Laura Barisoni
- Jens Hansen
- Ravi Iyengar
- Avi Rosenberg
- All KPMP members.

HuBMAP / IN U.

- Bruce Herr II
- Katy Borner

OBO Foundry

- Alex Diehl
- Bill Duncan

Funding:

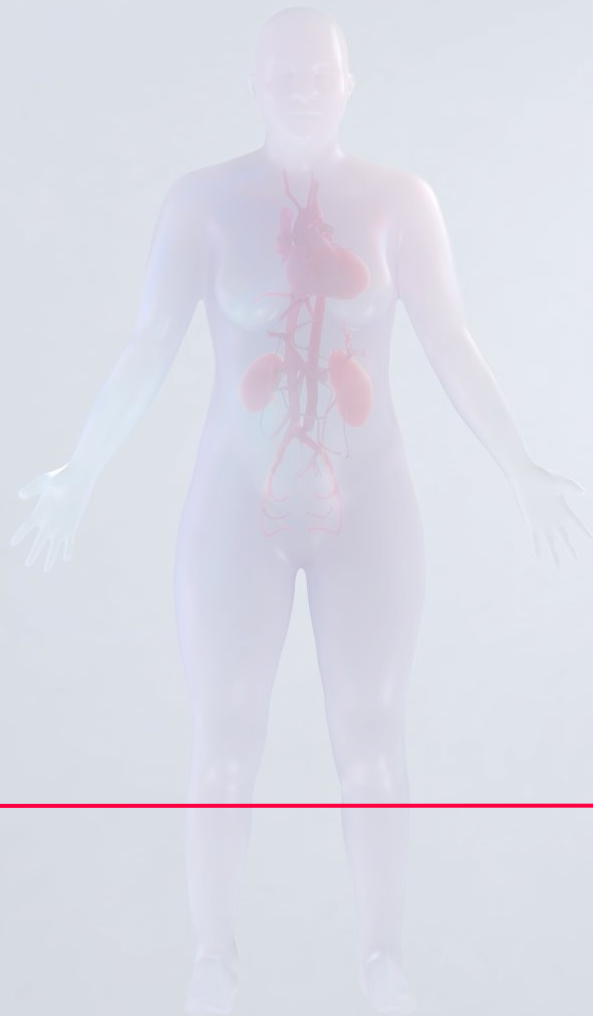
KPMP is funded from the NIDDK: U01DK133081, U01DK133091, U01DK133092, U01DK133093, U01DK133095, U01DK133097, U01DK114866, U01DK114908, U01DK133090, U01DK133113, U01DK133766, U01DK133768, U01DK114907, U01DK114920, U01DK114923, U01DK114933, U24DK114886, UH3DK114926, UH3DK114861, UH3DK114915, UH3DK114937.

Human Reference Atlas (HRA) research and development is funded by NIH OT2OD033756 and OT2OD026671, U24CA268108, U24DK135157 and U01DK133090.

The work has also been supported by the **Kidney Precision Medicine Project grant U2CDK114886**, HHSN316201300006W/HHSN27200002, and **HuBMAP U54 DK134301**.

NIAID grants to YH: R01AI081062; UH2AI132931; U24AI171008.

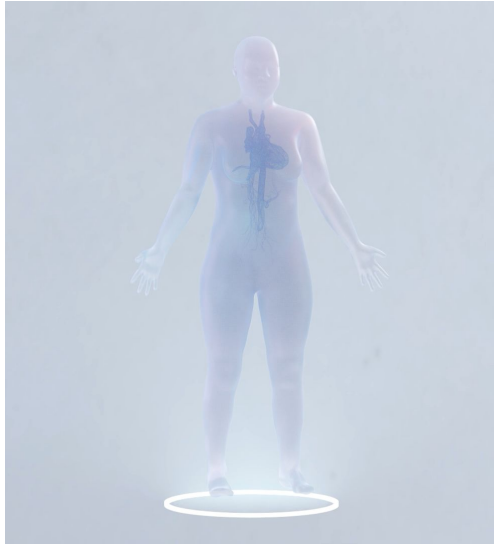
Q&A



<https://humanatlas.io/events/2024-24h>



Q&A



- How can ontologies improve AI, and how can AI improve ontologies?
- How can retrieval-augmented generation (RAG) specifically help with improving ontologies?
- How can AI chatbots (text-based) accurately represent ontologies (graph-based)?

Thank you
