

# SPOKE: Mapping Biomedical Knowledge for Precision Medicine and Genomics

<sup>1</sup>John H Morris, <sup>1</sup>Wanjun Gu, <sup>1</sup>Karthik Soman, <sup>1</sup>Rabia E Akbas, <sup>2</sup>Brett Smith, <sup>1</sup>Gabriel Cerono, <sup>1</sup>Braian Peetoom, <sup>1</sup>Catalina Villouta-Reyes, <sup>1</sup>Gundolf Schenk, <sup>1</sup>Angela Rizk-Jackson, <sup>3</sup>Lauren Sanders, <sup>3</sup>Sylvain V Costes, <sup>1</sup>Yongmei Shi, <sup>1</sup>Sharat Israni, <sup>2</sup>Sui Huang, <sup>4</sup>Peter W Rose, <sup>5</sup>Charlotte A Nelson, <sup>1</sup>Sergio E Baranzini

<sup>1</sup>University of California, San Francisco <sup>2</sup>Institute of Systems Biology <sup>3</sup>NASA <sup>4</sup>University of California, San Diego <sup>5</sup>Matebio.ai

UCSF

Bakar Computational Health Sciences Institute

## Introduction

Vast amounts of biomedical data are being generated at an unprecedented rate, leading to rapid growth in databases and repositories. Biomedical information is fragmented across disciplines due to the complexity of human physiology, limiting accessibility and interpretation.

Knowledge graphs are gaining traction as a method to integrate disparate biomedical data sources and facilitate knowledge discovery. The biomedical field requires a structured approach to connect diverse datasets and overcome compartmentalization. Connecting relevant biomedical information could lead to new insights in disease understanding, drug discovery, and personal health.

To this end, we introduce **SPOKE (Scalable Precision Medicine Open Knowledge Engine)**: A heterogeneous biomedical knowledge graph linking molecular biology, pharmacology, and clinical practice. By connecting diverse biomedical information, SPOKE aims to facilitate knowledge emergence, enhance disease understanding, improve drug discovery, and support personalized healthcare.

## Method

| Source                                | Node Type   |
|---------------------------------------|---|
| Cell Ontology                         | CellType  |
| Disease Ontology                      | Disease   |
| DrugBank, ChEMBL, KEGG                | Compound  |
| DrugCentral                           | PharmacologicClass                                      |
| Entrez Gene                           | Gene  |
| Enzyme Commission                     | EnzymaticActivity                                       |
| FoodOn                                | Food  |
| Gene Ontology                         | BiologicalProcess, CellularComponent, MolecularFunction |
| Geonames, Location databases          | Location  |
| MeSH, HPO                             | Symptom   |
| MiRDB                                 | MiRNA   |
| NCBI Taxonomy, BV-BRC                 | Organism  |
| NHANES                                | Blend, DietarySupplement                                |
| Pathway databases                     | Reaction  |
| Pfam                                  | ProteinDomain, ProteinFamily                            |
| SIDER                                 | SideEffect  |
| Uberon                                | Anatomy   |
| UniProtKB                             | Protein   |
| WikiPathways, Reactome, KEGG, MetaCyc | Pathway   |

- SPOKE Aggregates data from **51 specialized databases** across multiple biomedical domains.

- Data is structured into **33 node types** (e.g., genes, proteins, diseases, pathways) and **90 edge types** (relationships between nodes).

- **Weekly** updates for most data types, with some static sources.

### Network Structure and Knowledge Graph Design

**Heterogeneous Network:** Nodes represent different data types; edges define known relationships.

**Hierarchical Organization:** Directed edges indicate **hierarchies** (e.g., disease subtypes, anatomical relationships).

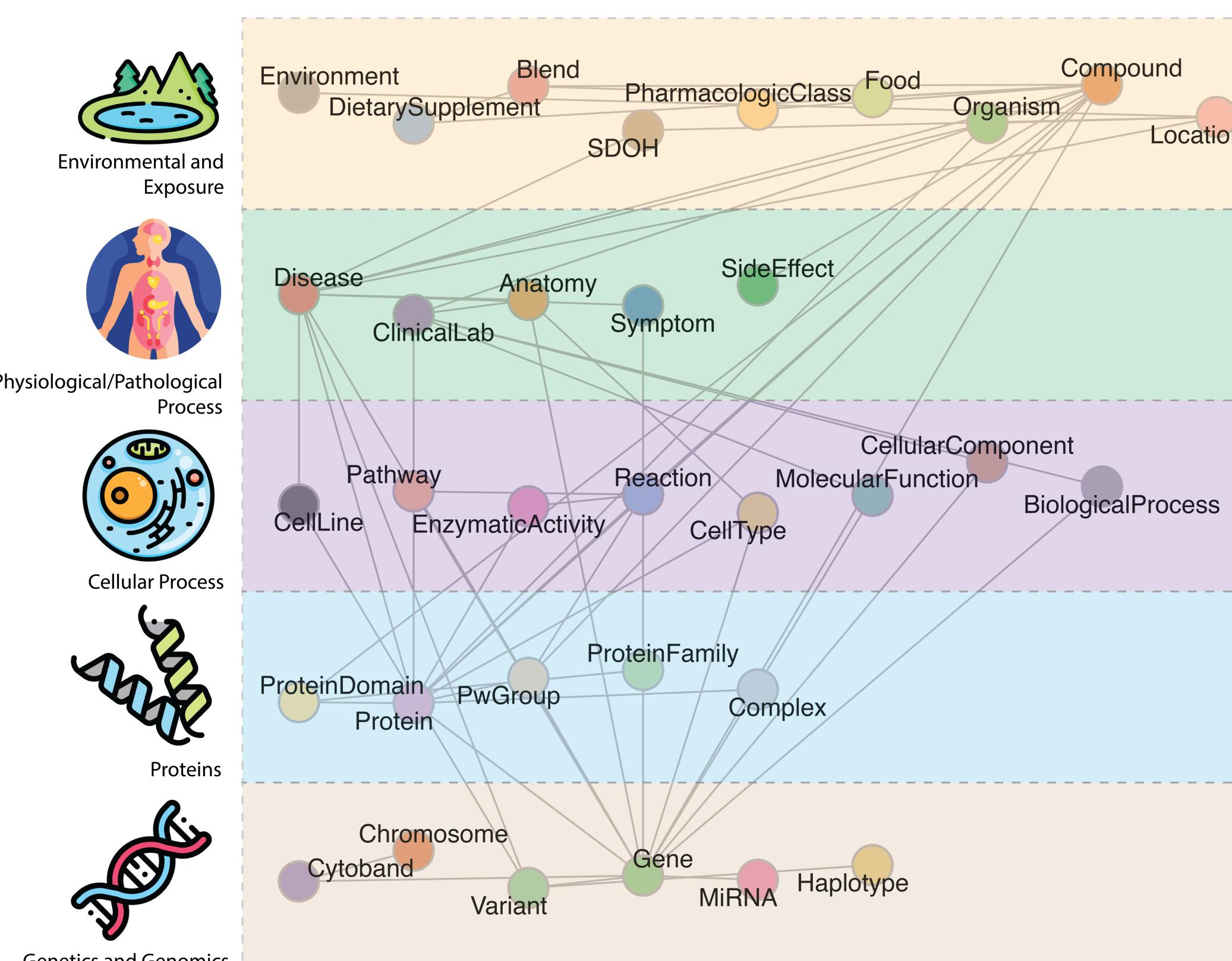
**Semantic and Functional Connectivity:** Identifies **unknown connections** by linking previously unrelated nodes through multi-step paths.

**Table 1 (left).** A representative list of data sources of all nodes in SPOKE

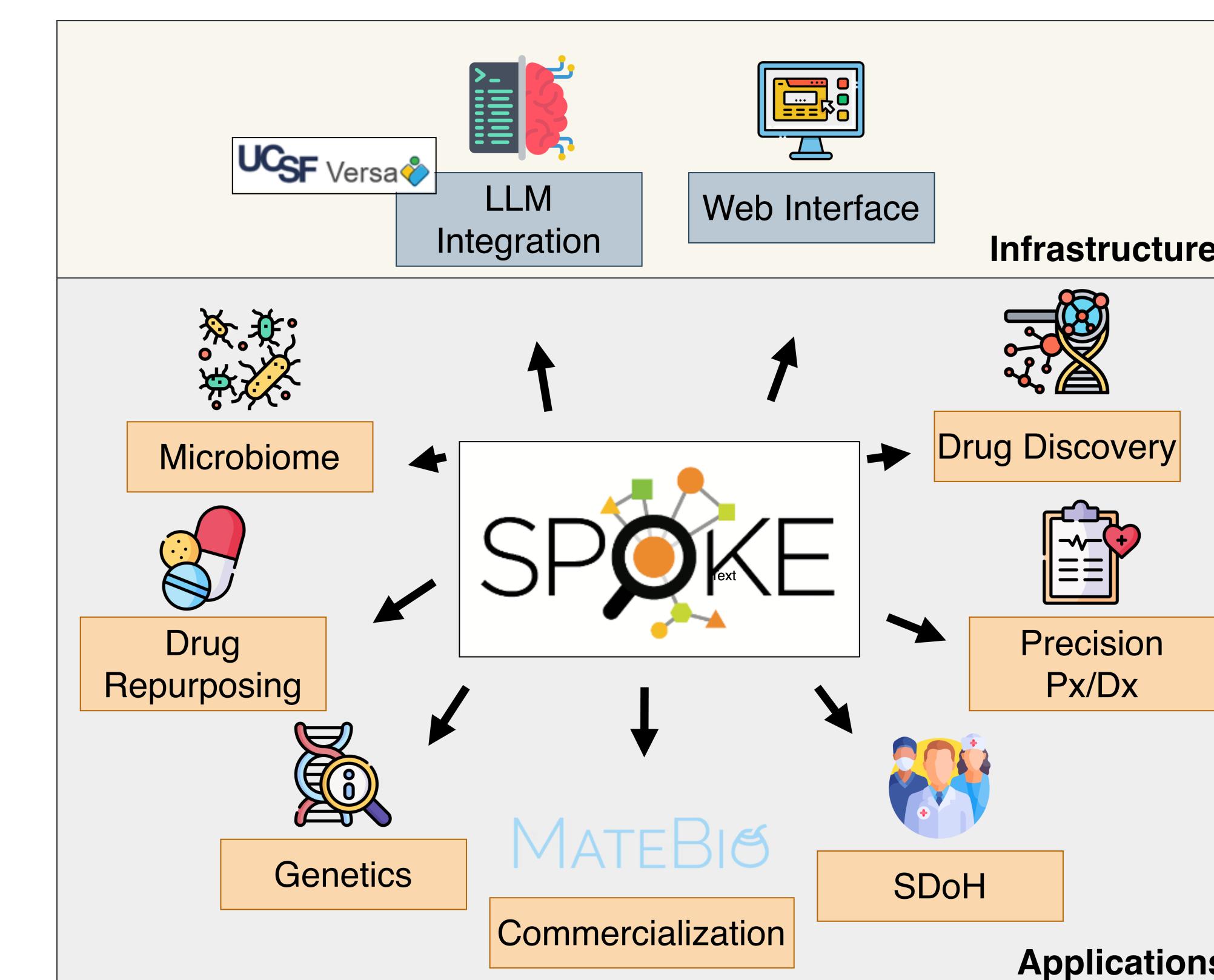
| Source                         | Edge Types                   |
|--------------------------------|------------------------------|
| Bgee                           | Anatomy-expresses-Gene       |
| BindingDB, DrugCentral, ChEMBL | Compound-binds-Protein       |
| BioGRID, STRING, IntAct        | Protein-interacts-Protein    |
| ChEMBL, DrugCentral            | Compound-treats-Disease      |
| ClinVar, GWAS Catalog          | Variant-associates-Phenotype |
| Gene Ontology, Pathway Commons | Gene-participates-Pathway    |
| MiRDB                          | MiRNA-targets-Gene           |
| OMIM, GWAS Catalog, DISEASES   | Disease-associates-Gene      |
| TFLink, TRRUST, ReMap          | Protein-regulates-Gene       |
| UniProt                        | Gene-encodes-Protein         |

**Table 2 (down).** A representative list of Data sources of all edges in SPOKE

## Infrastructure

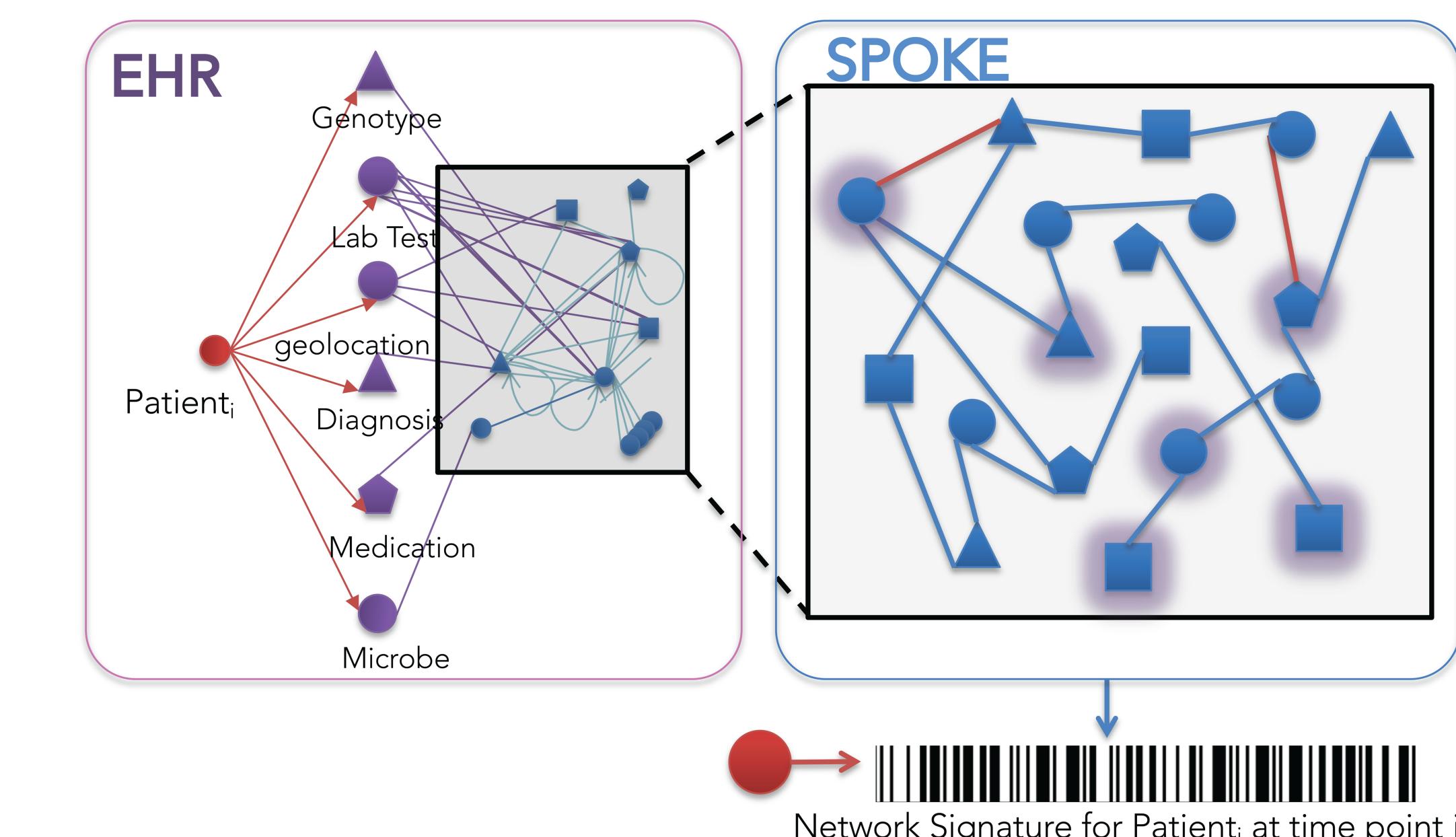


**Figure 1. Schematic of SPOKE.** Structure of the SPOKE knowledge graph, illustrating the interconnected biomedical entities across five hierarchical domains: **Environmental and Exposure**, **Physiological/Pathological Process**, **Cellular Process**, **Proteins**, and **Genetics/Genomics**. Nodes represent biological and medical concepts, while edges denote relationships between them, enabling integrative biomedical insights.

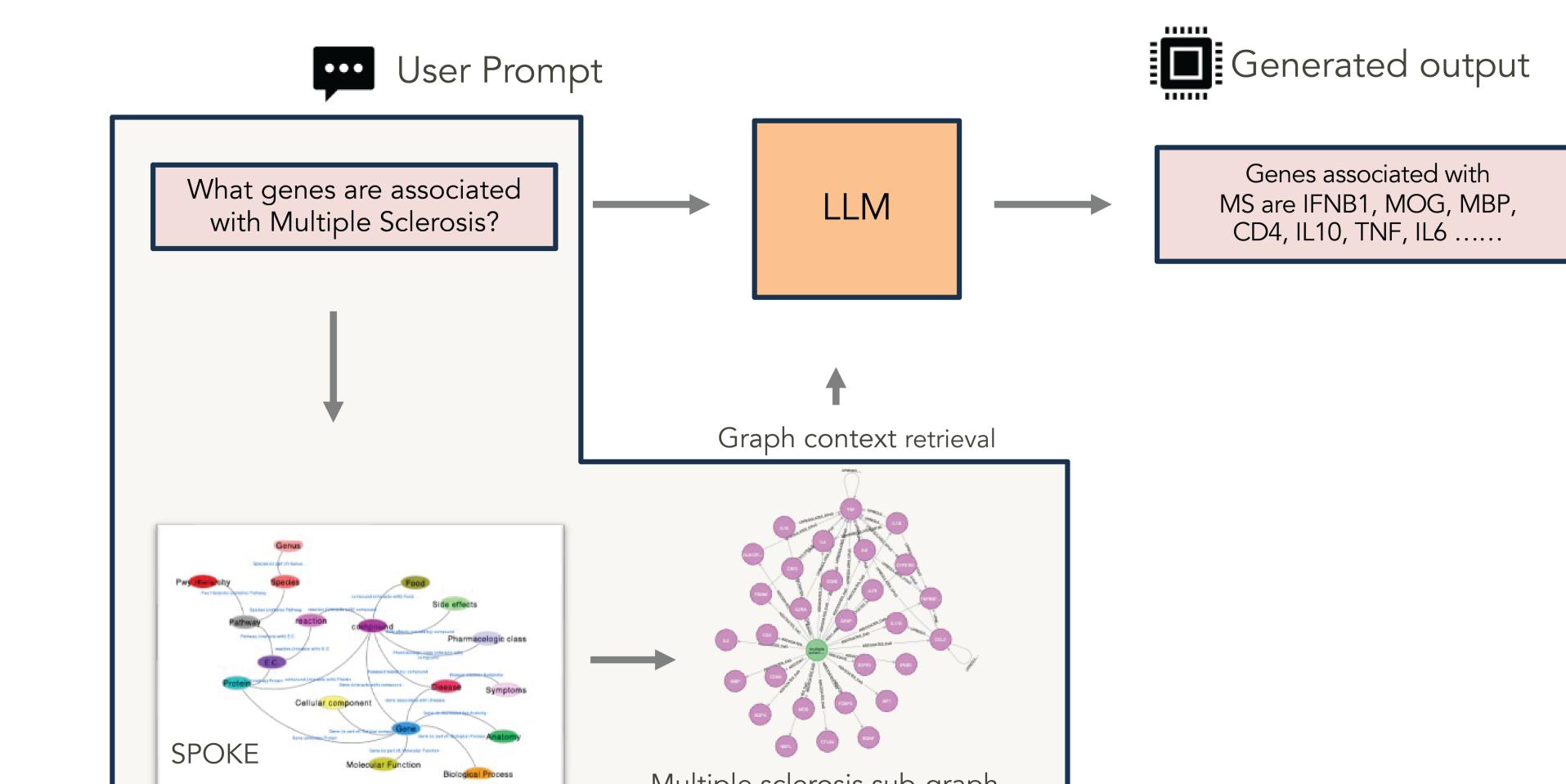


**Figure 2. Overview of SPOKE's infrastructure and applications.** The infrastructure integrates UCSF Versa with large language models (LLMs) and a web interface. Key applications include microbiome research, drug repurposing, genetics, social determinants of health (SDOH), precision medicine (Px/Dx), and drug discovery.

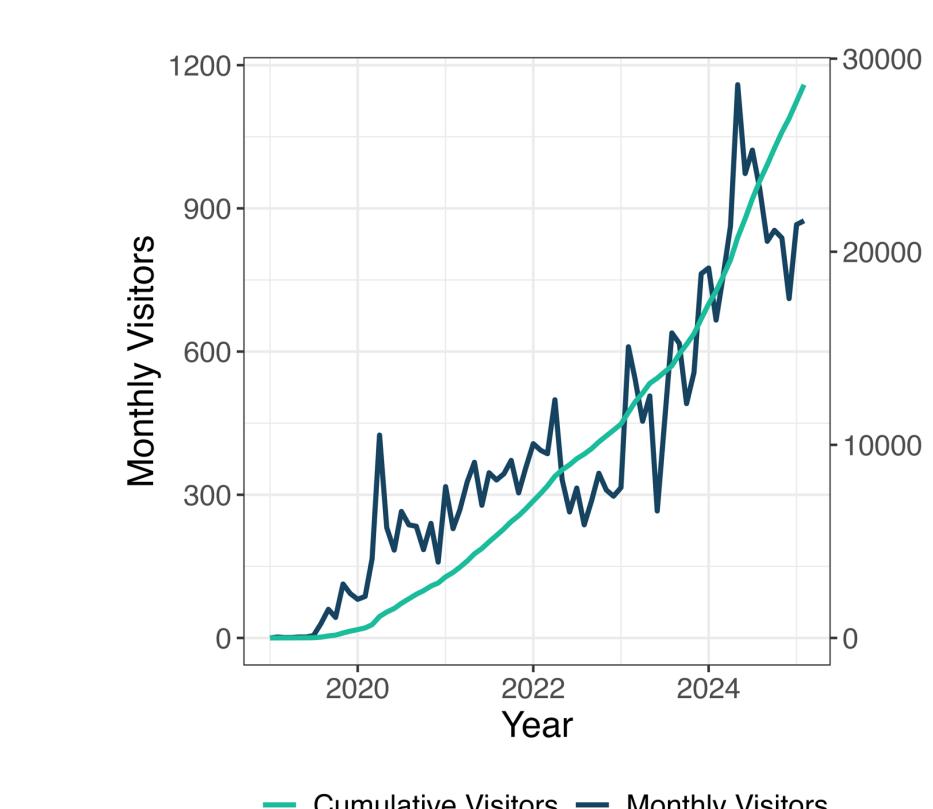
## Applications



**Figure 3. Integrating electronic health record (EHR) data into SPOKE for precision medicine.** Multi-modal patient data (e.g., genotype, lab tests, diagnoses, microbes) is embedded into the SPOKE knowledge graph, enabling network-based patient representations. These embeddings facilitate personalized disease progression prediction, subtyping, and diagnostics, supporting data-driven precision medicine.



**Figure 4. Integration of SPOKE with a large language model (LLM) using a retrieval-augmented generation (RAG) system.** When a user queries a biomedical concept (e.g., "What genes are associated with Multiple Sclerosis?"), the system retrieves relevant subgraphs from SPOKE, providing structured contextual information to the LLM. The LLM then synthesizes an informed response based on graph-derived knowledge, enabling accurate and interpretable answers to biological and medical inquiries.



**Figure 5. Growth in SPOKE's popularity over time.** The dark blue line represents the number of monthly visitors, showing an increasing trend with fluctuations. The green line represents cumulative visitors, indicating a steady rise in SPOKE's user engagement and adoption since its inception.

**Other Applications:** A subnetwork of SPOKE has been adapted into a pathway analysis tool, enabling the interpretation of omics data (transcriptomics, metabolomics, proteomics) to uncover molecular biology insights.