

# RARe-SOURCE™: Integrated Bioinformatics Resource for Rare Diseases

Jason Cheung<sup>1</sup>, Uma Mundunuri<sup>2</sup>, PhD, Daniel Watson<sup>2</sup>, Mohammad Alodadi<sup>2</sup>, PhD, Erica Lyons<sup>2</sup>, PhD, Anney Che<sup>2</sup>, Gregory Tawa, PhD<sup>1</sup>, London Toney<sup>1</sup>, Cara Purdy<sup>1</sup>, Richa Lomash<sup>1</sup>, PhD, Stephanie Mounaud<sup>1</sup>, Forbes Porter<sup>3</sup>, MD, PhD, Sharie Haugabook<sup>1</sup>, PhD, Elizabeth Ottinger<sup>1</sup>, PhD.

<sup>1</sup>Therapeutic Development Branch, Division of Preclinical Innovation, National Center for Advancing Translational Sciences, National Institutes of Health, 9800 Medical Center Drive, Rockville, MD, 20878

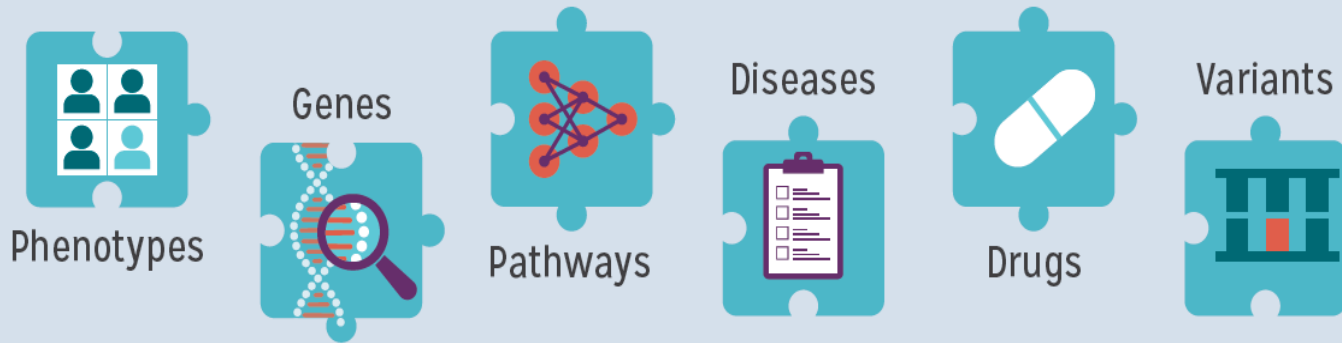
<sup>2</sup>Advanced Biomedical Computational Science, Frederick National Laboratory for Cancer Research, Frederick, MD, 21702

<sup>3</sup>Division of Translational Research, Eunice Kennedy Shriver National Institute of Child Health and Human Development, National Institutes of Health, Bethesda, MD, 20892



## Introduction

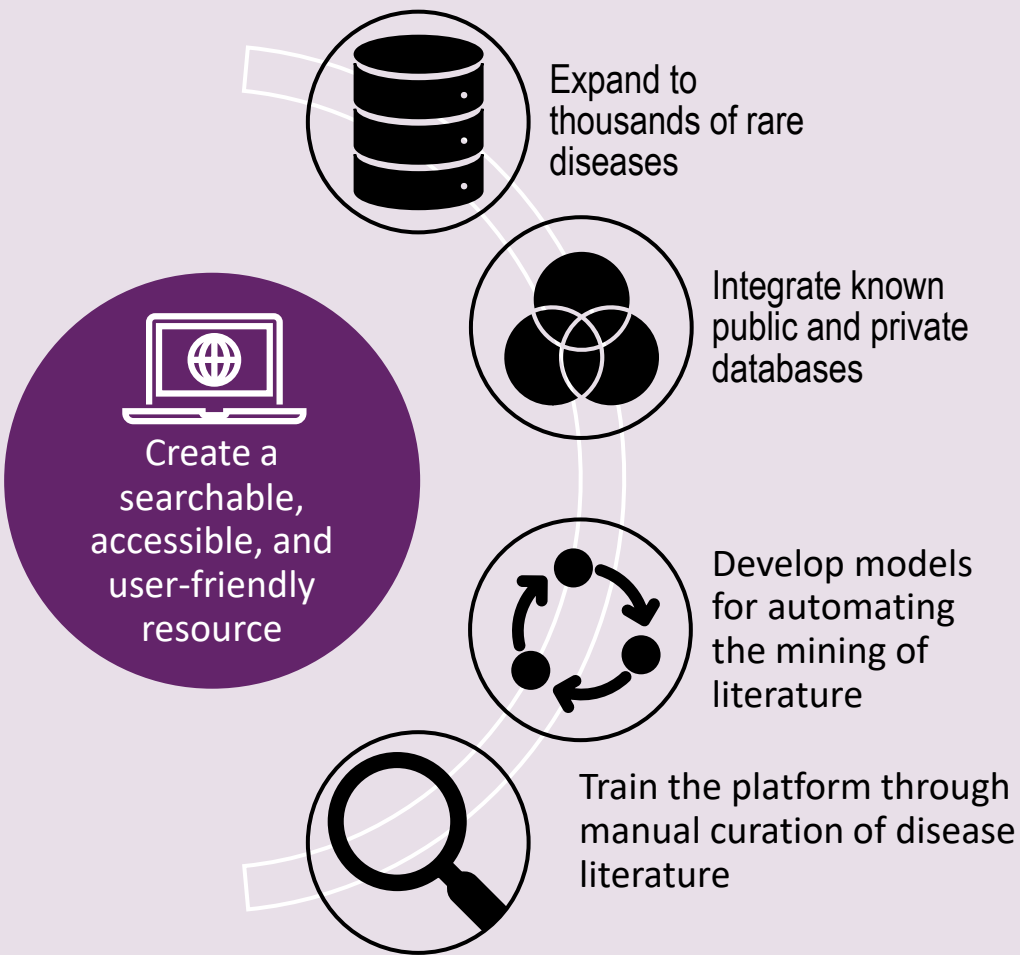
Rare disease data are fragmented and stored across different systems, platforms, and organizations. The ability to make connections between the diseases, their phenotypes, associated genes and related variants could unravel meaningful knowledge leading to the generation of novel hypotheses.



To harmonize and integrate these data sources, NCATS conceptualized RARe-SOURCE™, an integrated bioinformatics resource for rare diseases with the Advanced Biomedical Computational Science team at NCI-Frederick.



The approach is to extract, annotate and integrate disease-associated data from reputable sources including the peer-reviewed scientific literature to help end-users make molecular associations and use available data to catalyze the discovery and development of treatments for rare diseases.



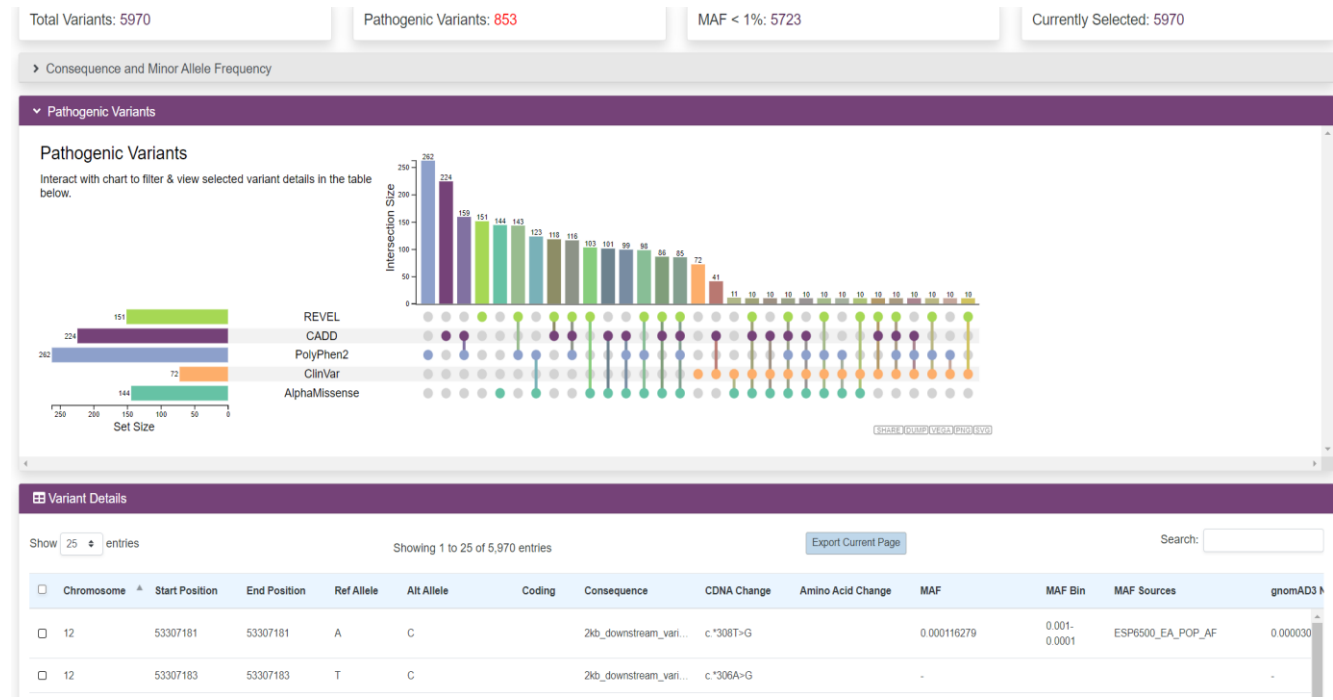
## Rare Disease / Gene Information

RARe-SOURCE™ integrates multiple data sources to provide comprehensive information on rare diseases with known genetic associations. Disease ontologies were integrated from GARD<sup>1</sup> and Orphanet<sup>2</sup>.

Rare Disease Name	Associated Genes	Gene Disease Annotations	Disease Annotations	Disease IDs	Links
X-linked Creatine Transporter Deficiency • Creat1 • Cerebral Creatine Deficiency Syndrome 1	SLC6A8	Curated Variants (186) NB Gene Disease Literature AI	NB Disease Literature AI	GARD: 0001608 OMIM: 300032 Orphanet: 520201	NB Published
12q14 Microdeletion Syndrome • 12q14 Microdeletion Syndrome • Del(12)(q14)	HMG2A LEMD3	NB Gene Disease Literature AI	NB Disease Literature AI	GARD: 0013390 Orphanet: 94963	NB Published
15q11.2 Microdeletion Syndrome • 15q11.2 Duplication Syndrome • 15q11.2 Microdeletion	NIPA1 TUBG1 NIPA2	NB Gene Disease Literature AI	NB Disease Literature AI	GARD: 0010025 OMIM: 615605 Orphanet: 261183	NB Published
15q13.3 Microdeletion Syndrome • 15q13.3 Microdeletion Syndrome	CHRNA7	NB Gene Disease Literature AI	NB Disease Literature AI	GARD: 0010296 OMIM: 612001 Orphanet: 190319	NB Published
15q24 Microdeletion Syndrome • 15q24 Deletion • 15q24 Microdeletion Syndrome	SRD5A	NB Gene Disease Literature AI	NB Disease Literature AI	GARD: 0012219 OMIM: 613405 Orphanet: 94065	NB Published
16p24.3 Microdeletion Syndrome • 16p24.3 Microdeletion Syndrome • Chromosome 16p24.3 Microdeletion Syndrome	ANKRD11	NB Gene Disease Literature AI	NB Disease Literature AI	GARD: 0010035 Orphanet: 261250	NB Published
17p11.2 Microduplication Syndrome • 17p11.2 Duplication Syndrome • 17p11.2 Microduplication Syndrome	RAI1	NB Gene Disease Literature AI	NB Disease Literature AI	GARD: 0010145 OMIM: 610803 Orphanet: 1713	NB Published

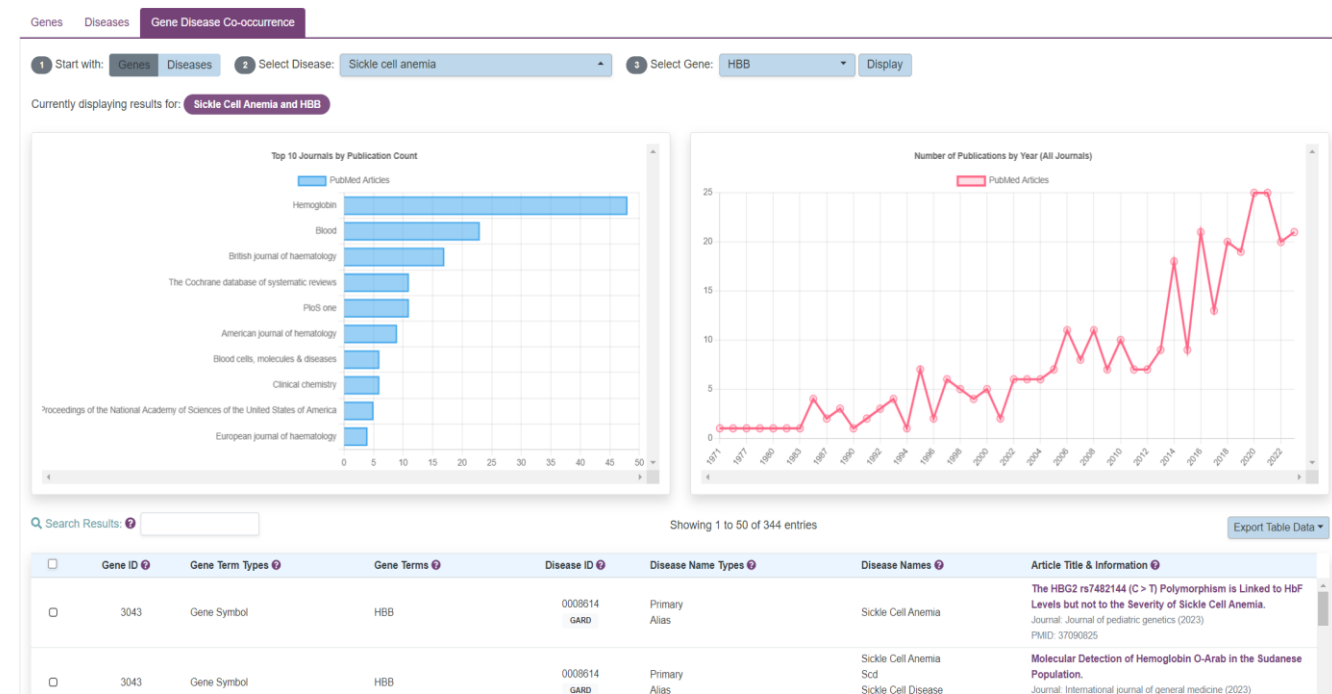
## Gene Variant Annotation

Module includes millions of variants from public data sources such as gnomAD and ClinVar, annotated using OpenCRAVAT, for rare disease associated genes.



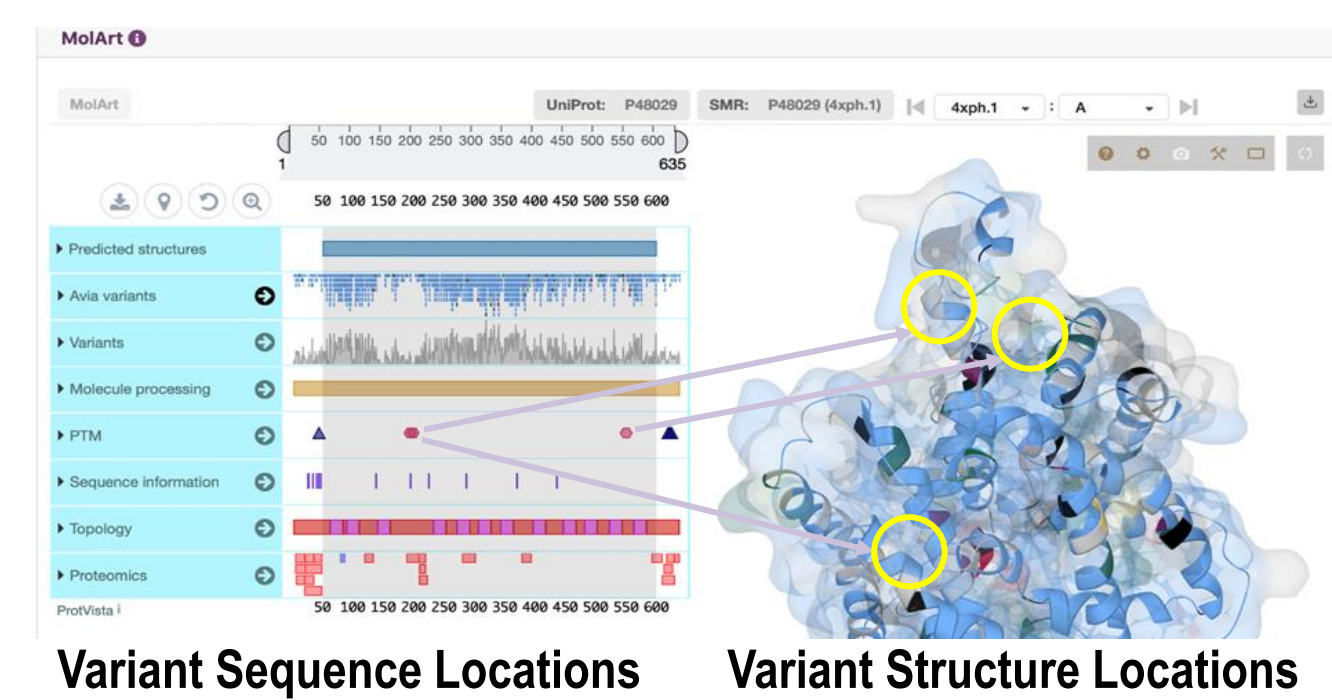
## Literature AI

Generates results from mining rare disease and associated gene mentions in all MEDLINE™ by using transformer models. It has an added benefit of easing the search process for finding literature on diseases or genes of interest.



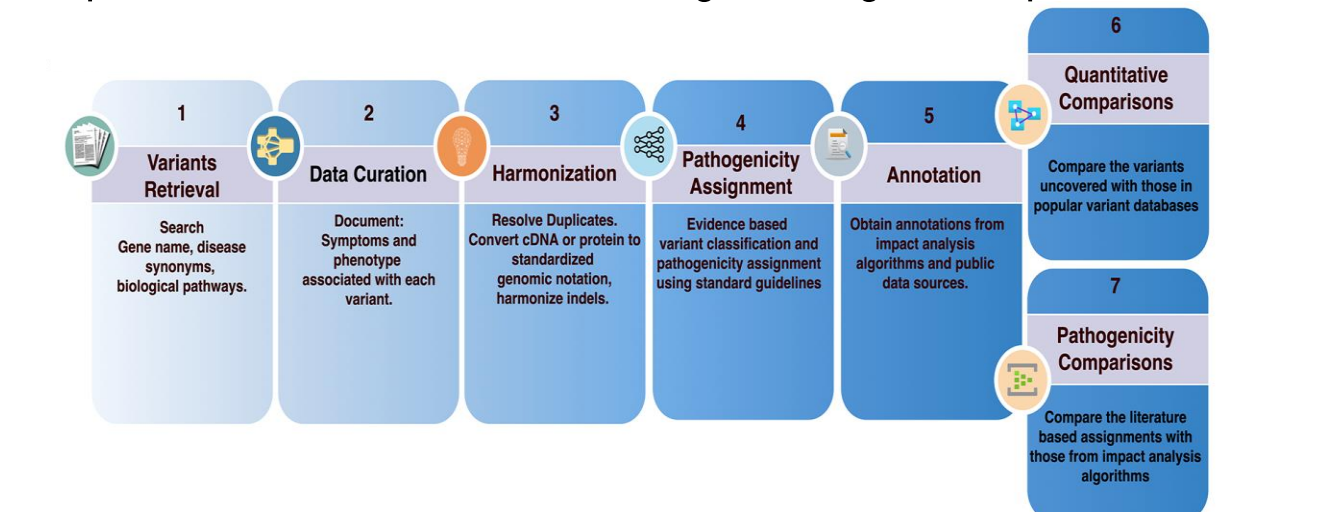
## 2D/3D Protein Structures

UniProt, MolArt, and Annotated Variants are integrated to illustrate the variant positions on the protein sequence and structure that can aid in functional impact analysis.



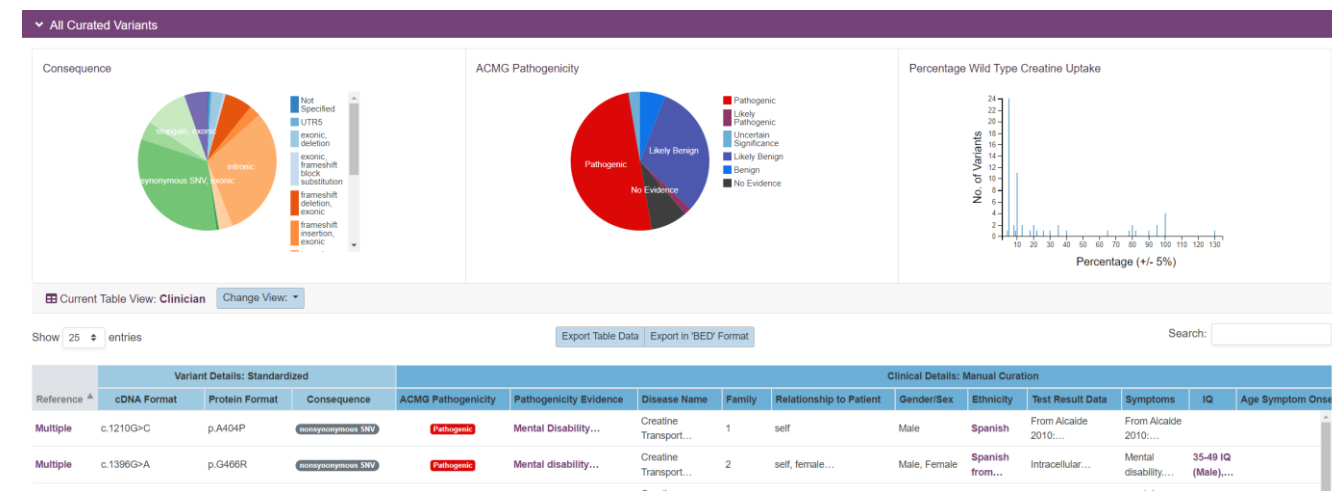
## Manual Curated Variants

The manual curation analyses of literature related to rare diseases serve as a vital validation tool for our Literature AI language model and provide a benchmark for thorough testing and improvement.



This involves extensive comparison of current and future AI language models with manual curation and tools like PubMed and PubTator.

Manual curation of SLC6A8 for creatine transporter deficiency<sup>3</sup> provides a curated list of published variants.

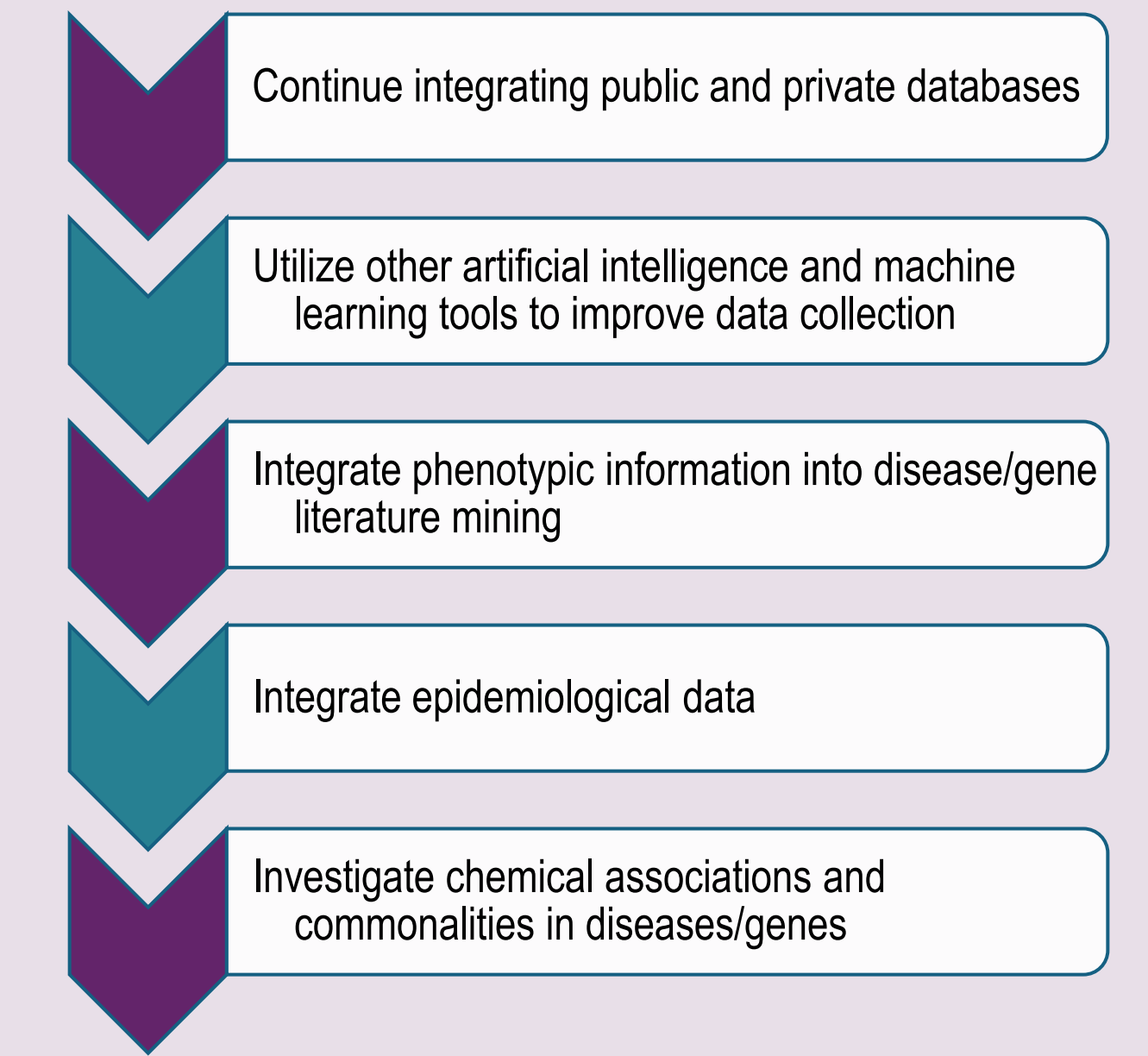


The curated details are reported as a highly annotated dataset of variants with clinical context, functional details and interactive visualizations.

## Future Directions

RARe-SOURCE™ is continuing to evolve by adding and integrating data and new features. The platform serves as a centralized information hub for rare diseases, to unlock novel insights into commonalities to identify new therapies to treat rare disorders.

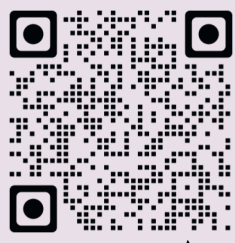
Plans include integrating genotype-phenotype mapping, pathway associations, responses to chemical compounds, and disease correlations.



## Contact

<https://raresource.nih.gov>

[rare-source@mail.nih.gov](mailto:rare-source@mail.nih.gov)



Scan Me!

**Notice:** Feedback or reports of inaccuracies within the resource are welcomed to facilitate data accuracy and inclusivity.

### References:

- [1] Genetic and Rare Diseases Information Center (GARD). <https://rarediseases.info.nih.gov/>
- [2] Orphanet: an online rare disease and orphan drug data base. <https://www.orpha.net>
- [3] Lyons, E.L., Watson, D., Alodadi, M.S. et al. Rare disease variant curation from literature: assessing gaps with creatine transport deficiency in focus. BMC Genomics 24, 460 (2023). <https://doi.org/10.1186/s12864-023-09561-5>

**Acknowledgements:** The authors thank NCATS leadership for their support of this project.

**Funding:** This project is being supported by the Intramural Research Program of NCATS/NIH under Contract No. HHSN261201500003 through NCI. This project has been funded in part with Federal funds from the NCI/NIH/HHS under Contract No. 75N91019D00024. The content of this publication does not necessarily reflect the views or policies of the Department of Health and Human Services, nor does mention of the trade names, commercial products, organizations imply endorsement by the U.S. Government.