# ClinCluster: Aggregating Disease Terms in ClinVar

NCBI's ML/AI Codeathon

Feb 26 - Mar 1, 2024

**National Library of Medicine**
*National Center for Biotechnology Information*

# ClinVar



Hypercholesterolemia due to variants in the gene LDLR

Any inherited type of hypercholesterolemia

Problem statement: Diseases in ClinVar are very granular and result in many variant-disease records.

Can we use an ML/AI approach to aggregate disease terms in ClinVar to reduce the number of variant-disease records?

# ClinCluster: Use LLM to aggregate these similar disease terms



**Storing and serving related documents with openai-to-sqlite and embeddings**

I decide to upgrade the related articles feature on my TILs site. Previously I calculated these using full-text search, but I wanted to try out a new trick using OpenAI embeddings for document similarity instead.

My openai-to-sqlite CLI tool already provides a mechanism for calculating embeddings against text and storing them in a SQLite database.

I was going to add a command for calculating similarity based on those embeddings... and then I saw that Benoit Delbosc had opened a pull request implementing that feature already!

I took Benoit's work and expanded it. In particular, I added an option for saving the resulting calculations to a database table.

This meant I could find and then save related articles for my TILs by running the following:

wget https://s3.amazonaws.com/til.simonwillison.net/tils.db

This grabs the latest tils.db used to serve my TIL website.

openai-to-sqlite embeddings tils.db \
  --sql 'select path, title, topic, body from til'

This retrieves and stores embeddings from the OpenAI API for every row in my til table - embedding the title, topic and body columns concatenated together, then keying them against the path column (the primary key for that table).

The command output this:

Fetching embeddings  [###############################]  100%
Total tokens used: 402500

402,500 tokens at $0.0001 / 1K tokens comes to $0.04 - 4 cents!

Now that I've embedded everything, I can search for the most similar articles to a particular article like this:

```
[
    0.0017127031460404396,
    -0.0049754707142710690,
    0.0107359681127846718,
    -0.0079374928027391430,
    -0.0177794024348258970,
    0.0166018623858690260,
    -0.0177794024348258970,
    ...
    Array of floating point numbers
]
```
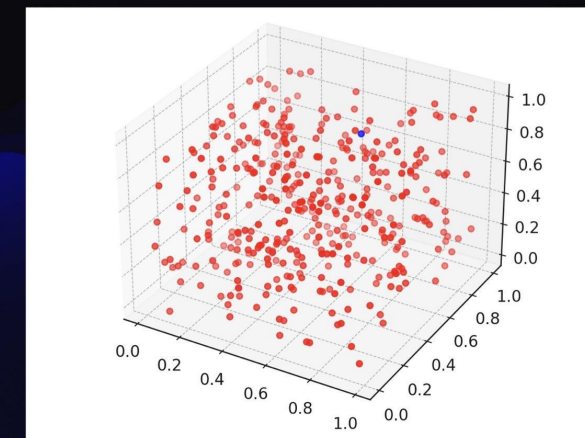Fixed size: 300, 1000, 1536…

**Embeddings in LLM**: take a piece of content and turn that piece of content into an array of floating point numbers.

A location in 1,536 dimension space



Same length, no matter how long the content is

**Workflow**
1. Feed the disease names to LLM
2. LLM gives a embedding (a array of floating point numbers) for each name
3. Use the embeddings to cluster the disease terms

National Library of Medicine
National Center for Biotechnology Information

Images from https://simonwillison.net/2023/Oct/23/embeddings/

# DBSCAN: cluster the disease names



DBSCAN

k-means

Image from https://github.com/NSHipster/DBSCAN

Result:

- We forked and modified llm-cluster package
  https://github.com/simonw/llm-cluster
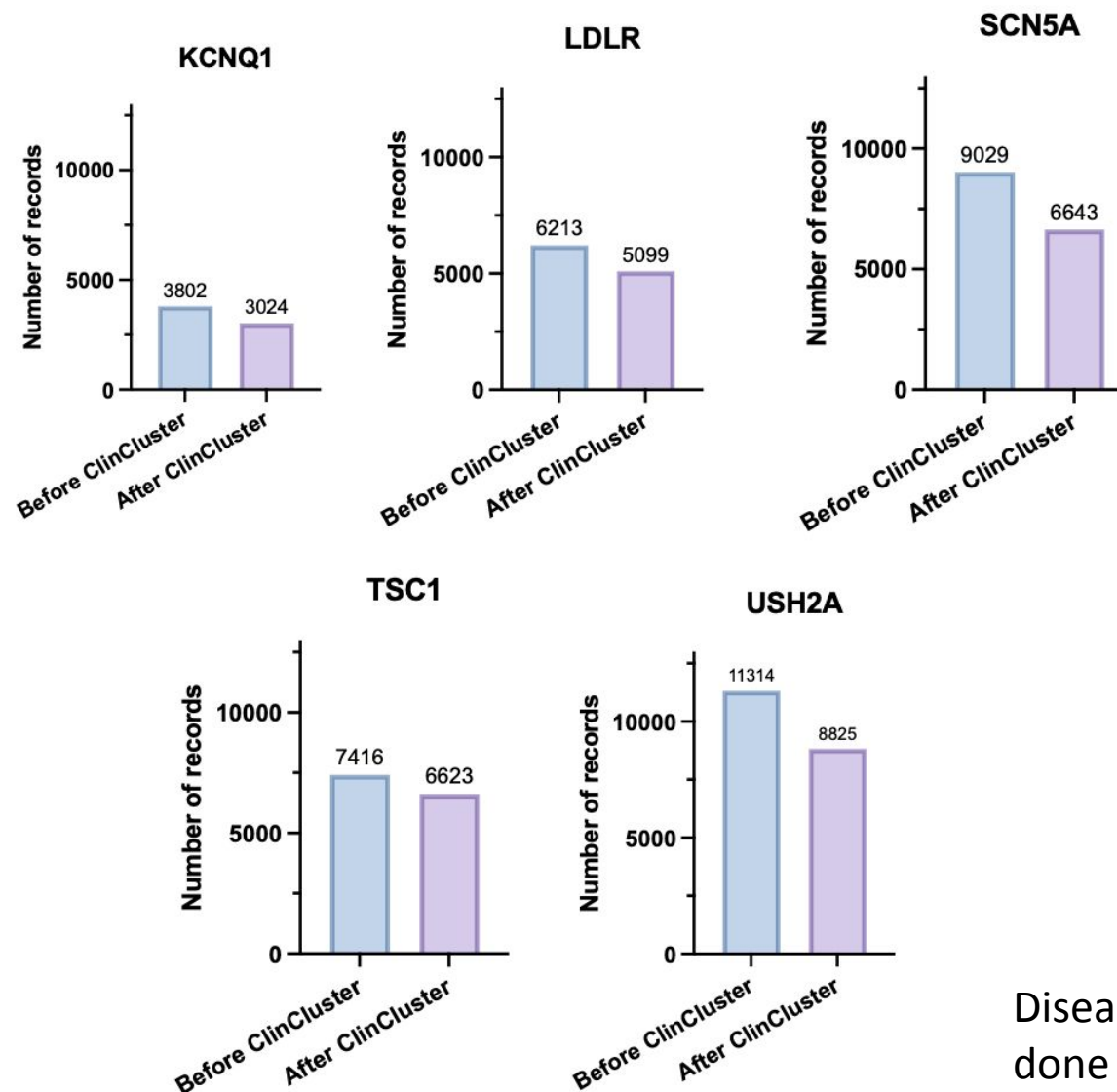- We change the k-means clustering algorithms to *unsupervised* DBSCAN algorithm for clustering the disease names.



NIH National Library of Medicine
*National Center for Biotechnology Information*
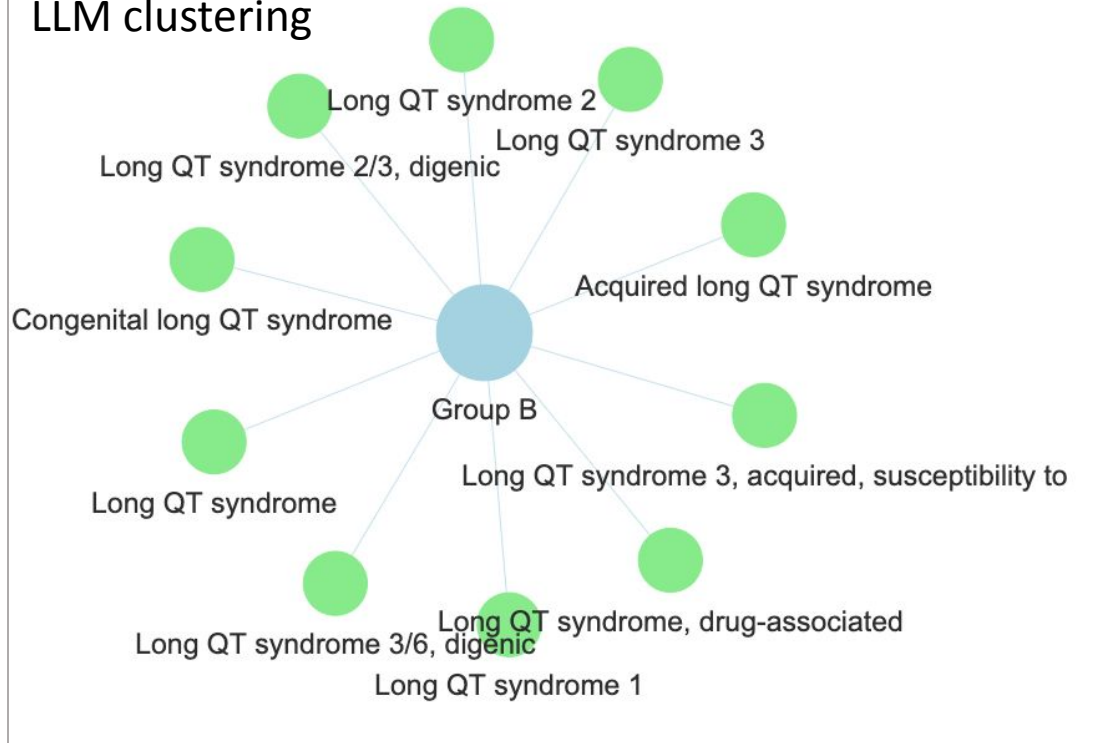
# ClinCluster: Performance metric



Significant reduction in the
total number of RCV records

Disease term clustering is
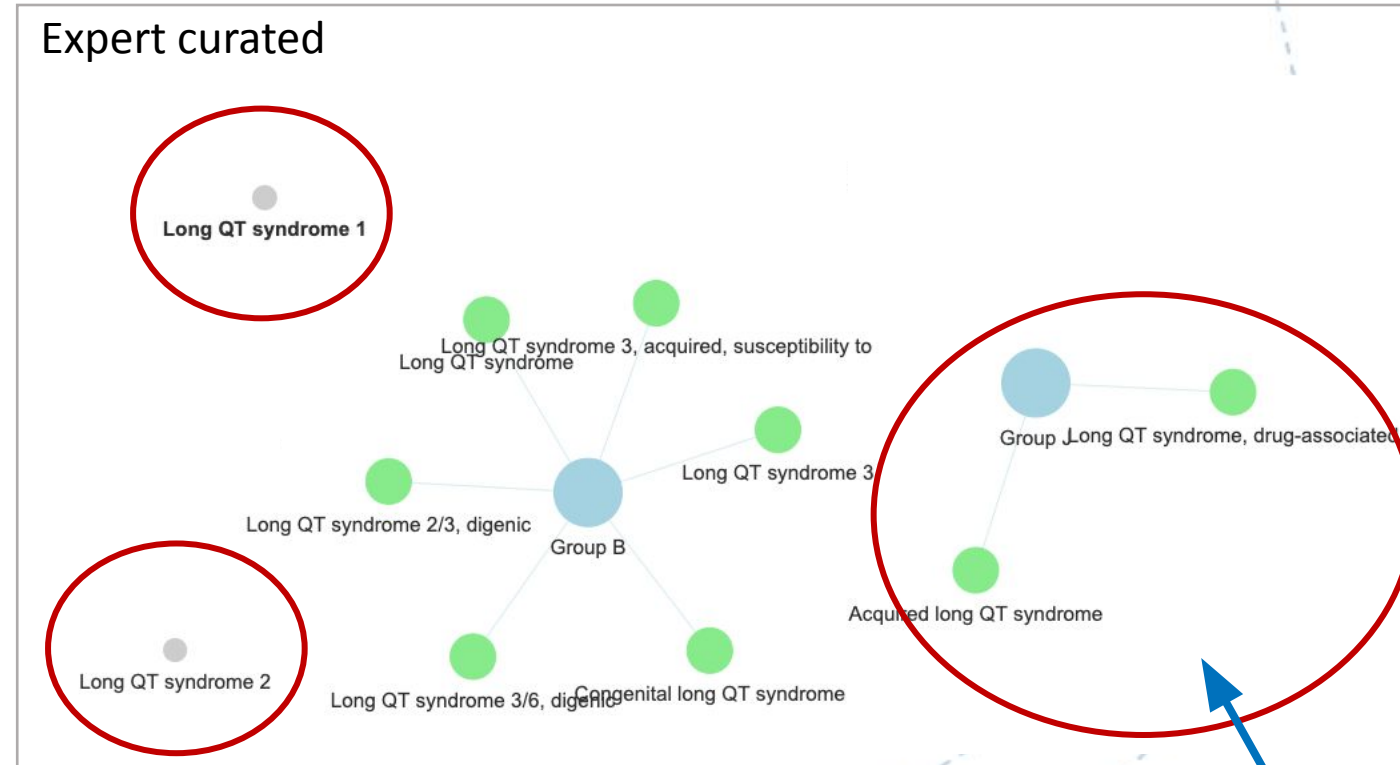done by LLM clustering algorithm.

# However, LLM clustering is not perfect…
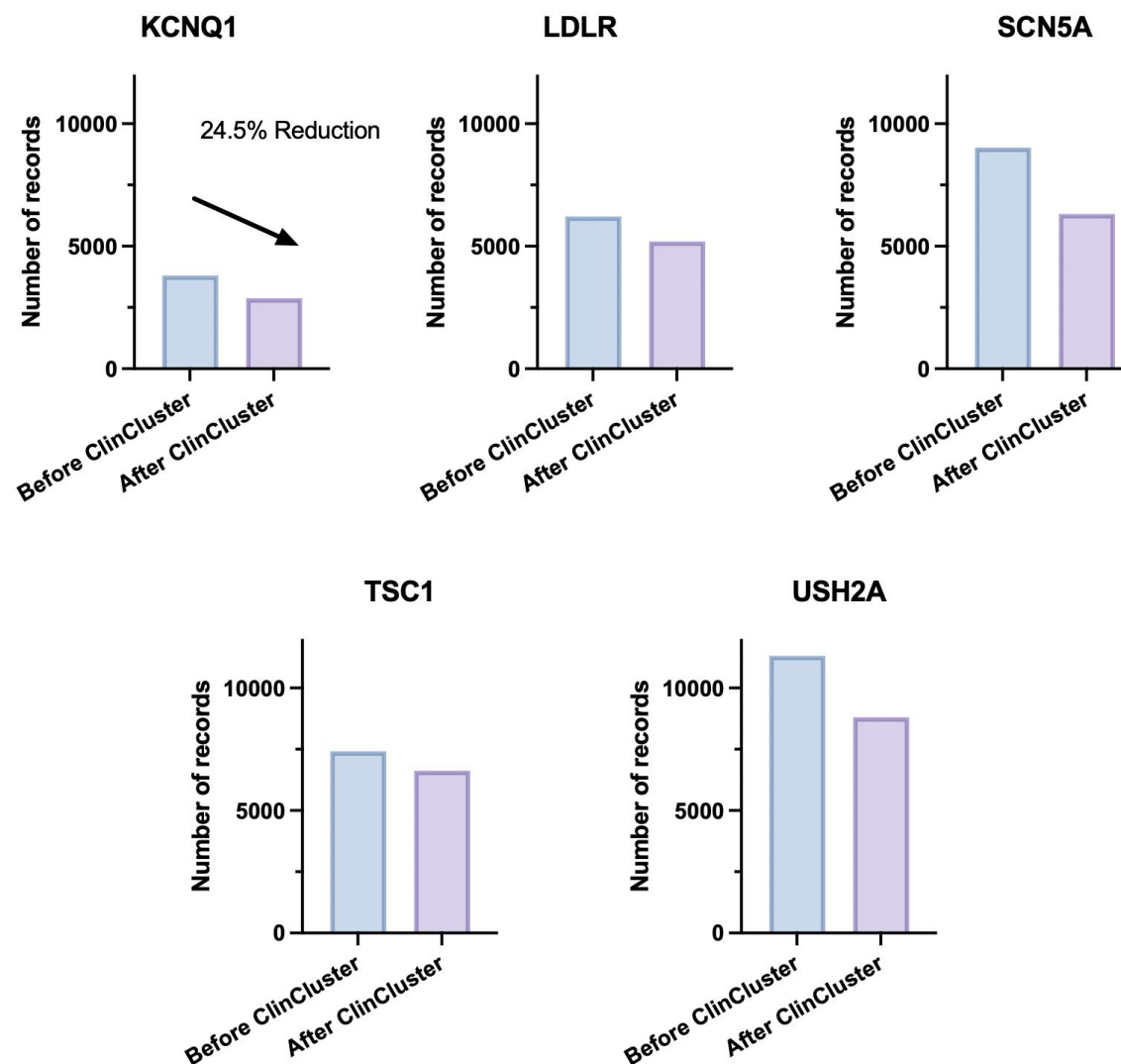
Disease clusters for SCN5A gene



Circled diseases are caused by other underlying conditions

State-of-art LLM might help

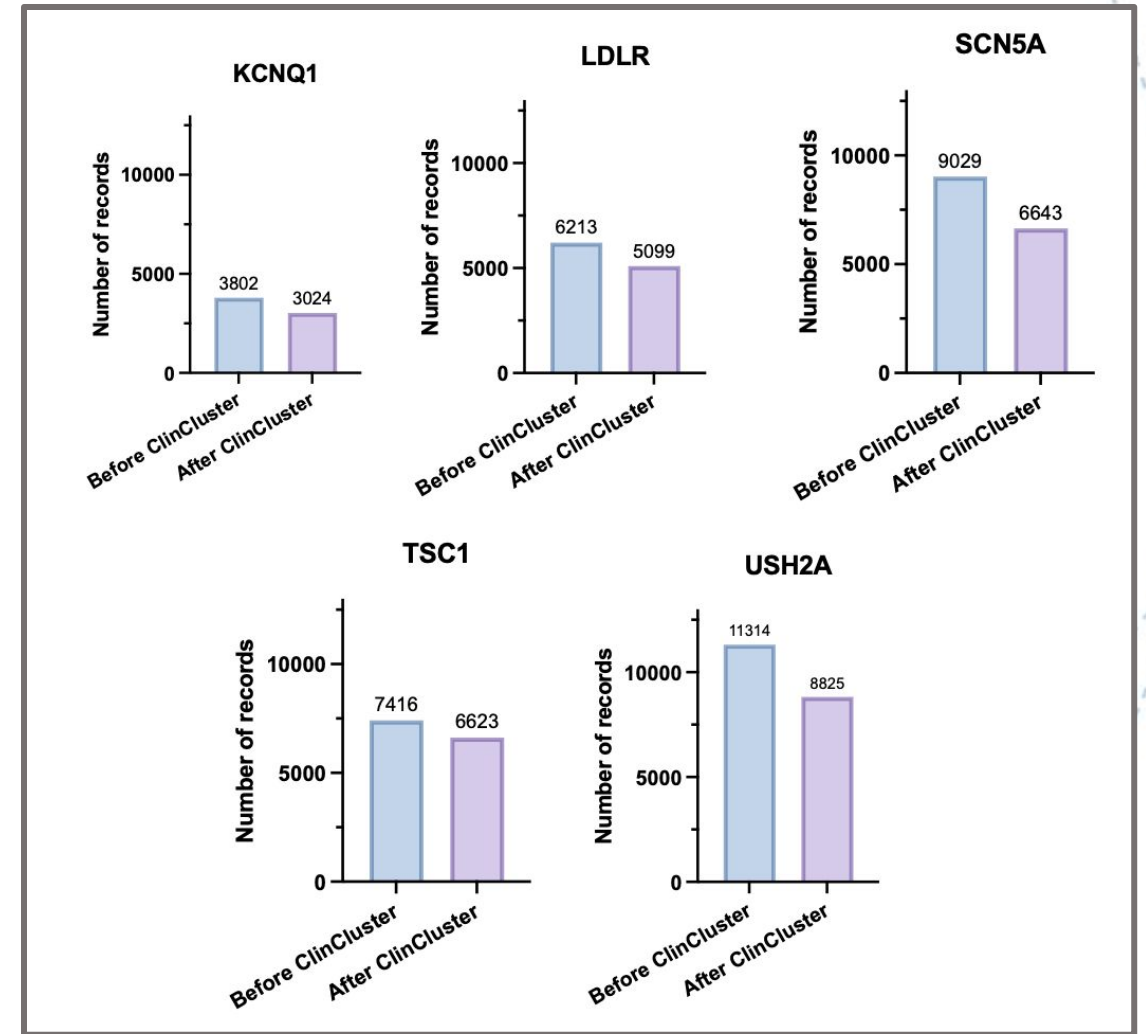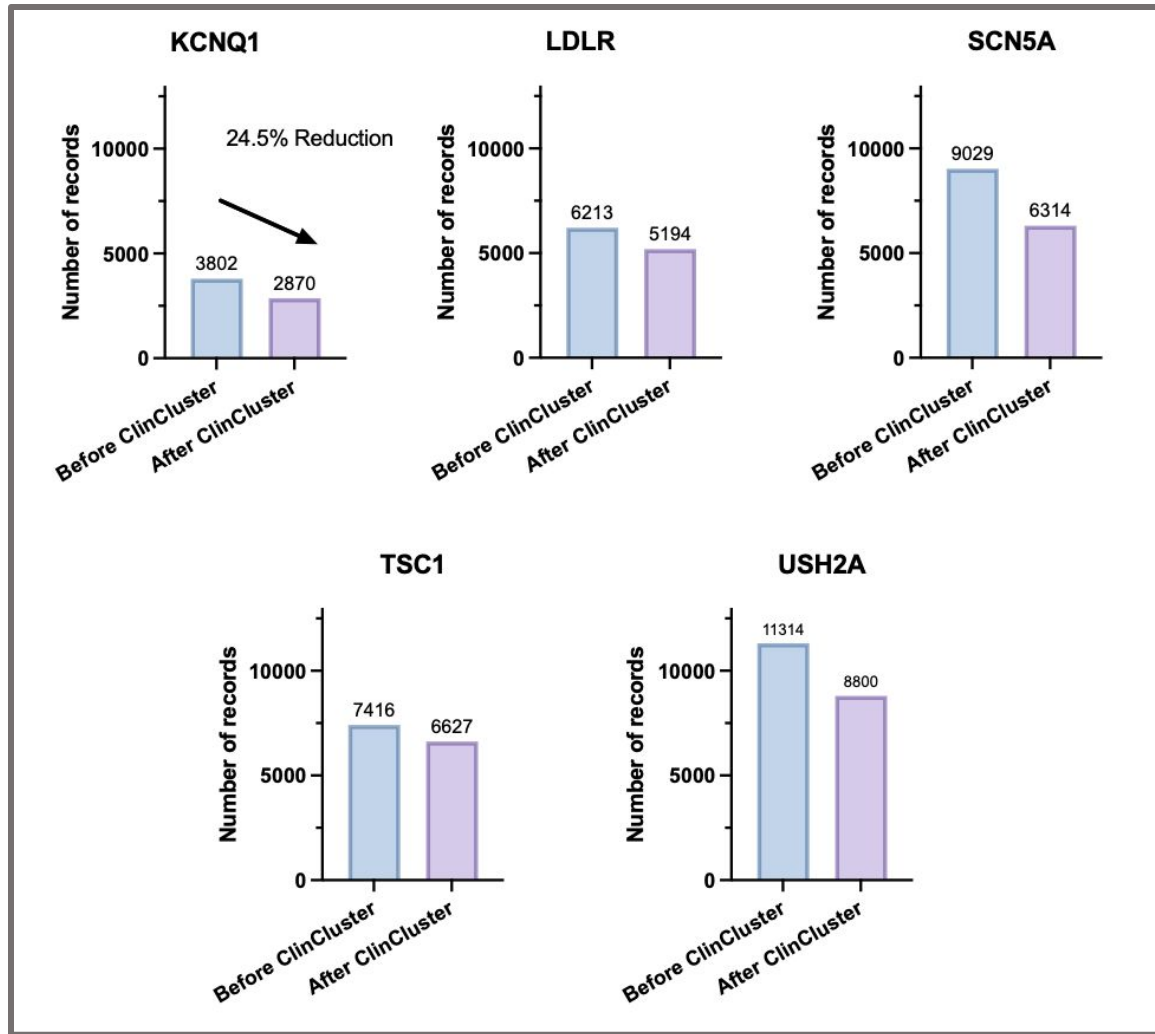Rare conditions

# ClinCluster: Performance metric



Significant reduction in the total number of RCV records

Disease term clustering is augmented by human expert

# ClinCluster: Performance metric



Disease term clustering is augmented by human expert after LLM clustering.
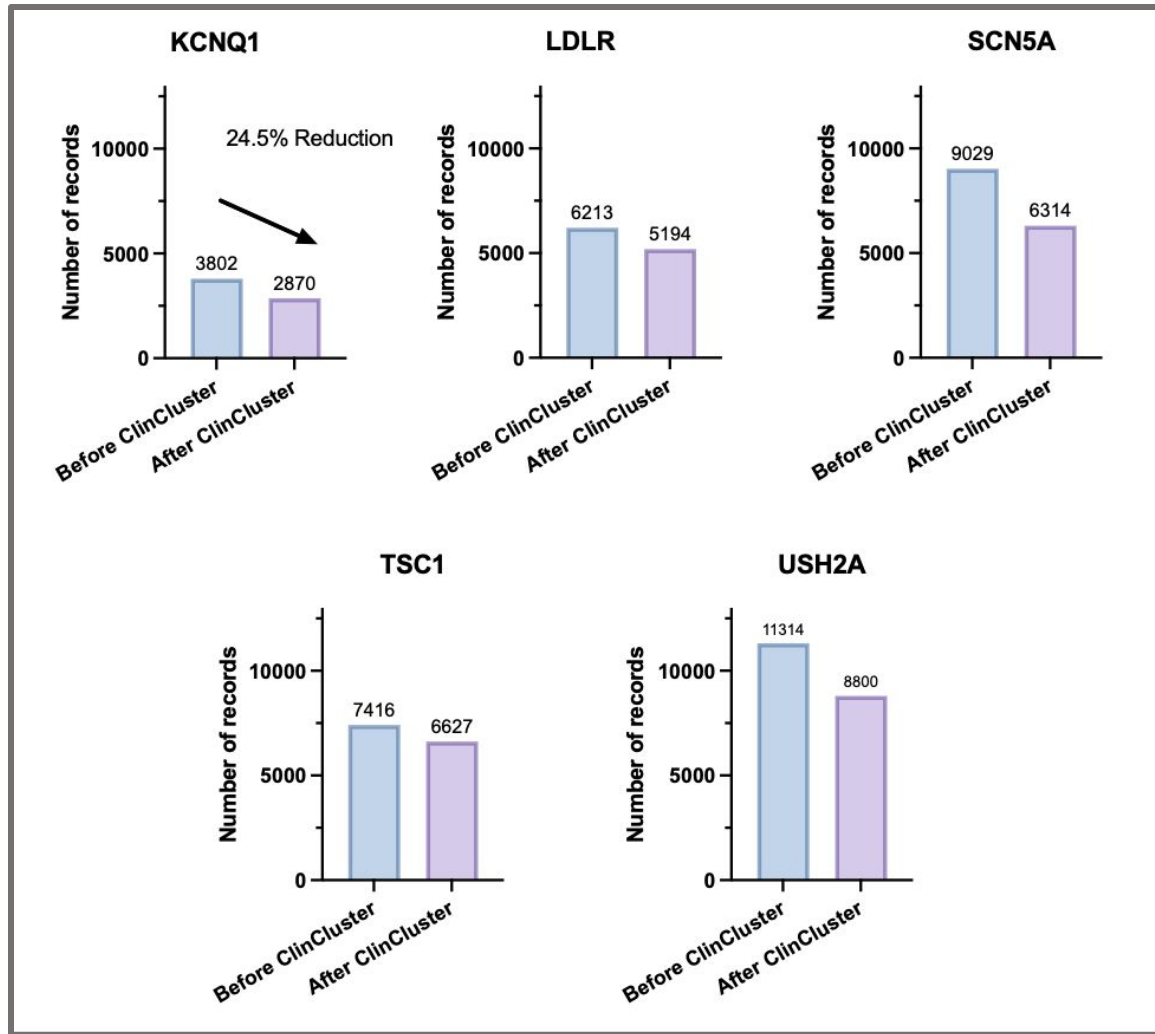
Disease term clustering is done by LLM clustering algorithm.

# ClinCluster: Performance metric

# Reproducible implementation

| Command Palette | ⌘ | P |
| Command Search | ^ | R |
| Warp AI | ^ | Space |

base ~

# Acknowledgements

Lauren Edgar, NIH/NHGRI

Benjamin Kesler, Vanderbilt University

Nicholas Minor, University of Wisconsin

Michael Muchow, Unaffiliated

Rebecca Orris, NIH/NCBI

Wengang Zhang, NIH/NCI

Guangfeng Song, NIH/NCBI (Co-Team Leader)

Melissa Landrum, NIH/NCBI (Co-Team Leader)

NIH〉 National Library of Medicine
National Center for Biotechnology Information