

Leveraging Multiple Public Unstructured Oncology Data Knowledge Compendium for Advanced Queries, Search and Predictions



Vishakha Sharma¹, Alex Thomas², Vishnu Vetrivel², David Talby², Antoaneta Vladimirova¹
¹Roche Diagnostics, Roche Information Solutions, Santa Clara, CA; ²John Snow Labs, Seattle, WA

Abstract

In healthcare product development and research, teams invest huge amounts of time to study through publications and other relevant resources. There is a need for a novel solution to efficiently and reliably extract information from multiple clinical resources, in addition to generating new insights which can only be achieved through structuring textual information and accessible intelligent synthesis across multiple relevant resources. We created a cloud-based solution where data from heterogeneous sources is structured, integrated and harmonized, and users can easily leverage the combined database to answer domain-specific questions and generate insights efficiently in a targeted way.

Materials and Methods

Knowledge graphs provide us the advantage to encode and leverage relationships in addition to concepts in the context of heterogeneous data. We leveraged graph and NLP AI techniques to build a domain-specific knowledge graph. We extracted the biomedically-relevant subset of wikidata, and augmented it by adding more entities and relationships from the biomedical literature (PubMed), clinical trials (clinicaltrials.gov) and NIH grants. We leveraged domain-specific named entity recognition (NER) models to identify and include rich biomedical entities.

The data is stored in data store (the graph itself), search indexes (the documents), and database tables (derived data for the visualizations).

We used an embedding model of terms and MeSH entities in order to create the scatter plot of related terms in the trend visualization. The trending terms are looking at year-over-year percentage increase in occurrences in the select set of documents.

The biomarker model is produced using features from a TransE-L2⁵ embedding and a classification model. In order to make the problem tractable, the possible pairs are limited to those connected by a fixed set of paths.

Results

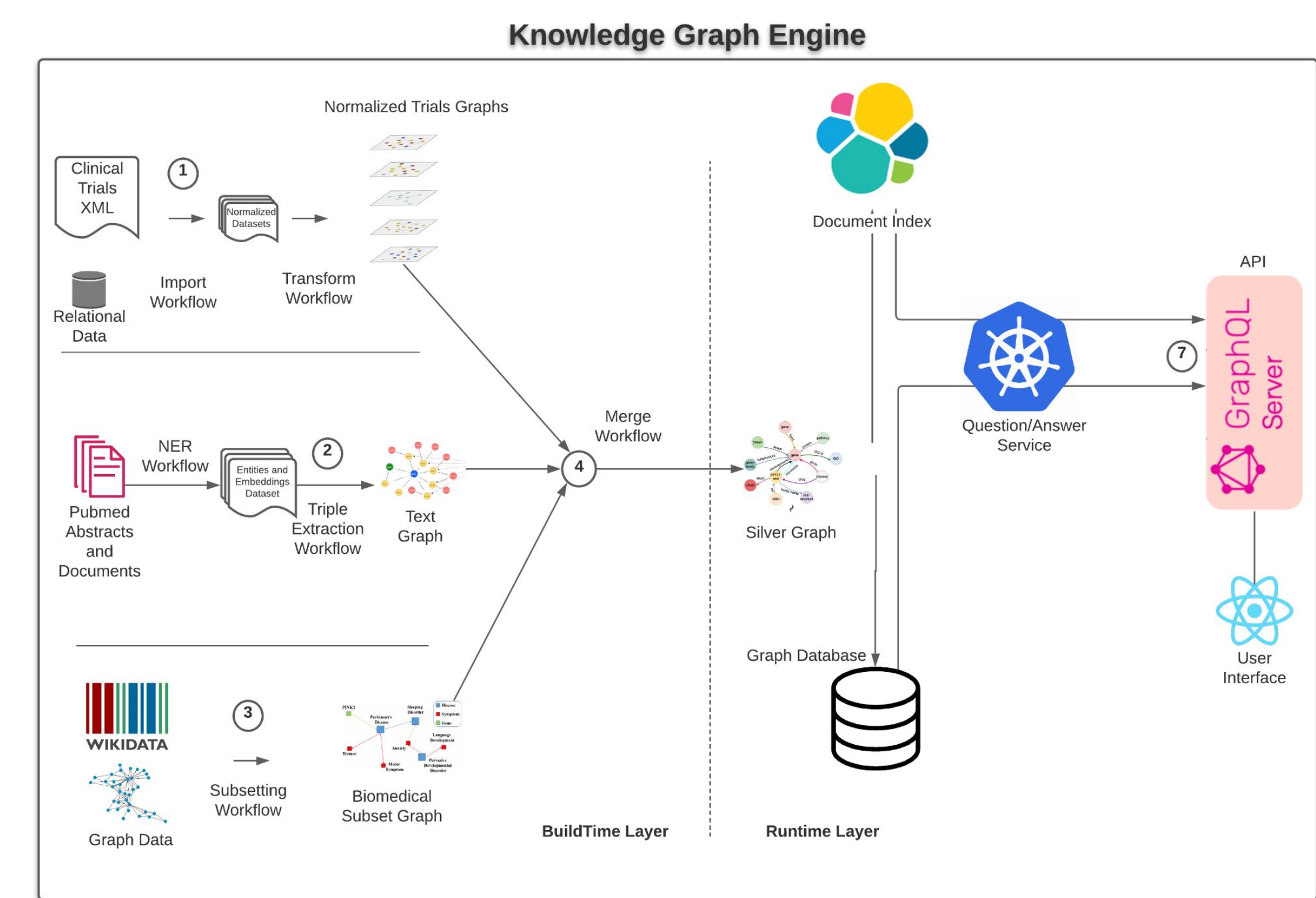


Figure 1. Schematic Diagram of Knowledge Graph Engine

Data is extracted from structured data, text, and graphs, and stored in a graph-like tabular format. The data is then merged by mapping references to an entity to a common identifier (usually the wikidata QID). The properties and relationships are mapped to predicates existing in wikidata. Once this mapping has been completed, the data is merged into a single graph. The documents (e.g. from pubmed) are stored in their own search index. The graph itself is loaded into a triple store.

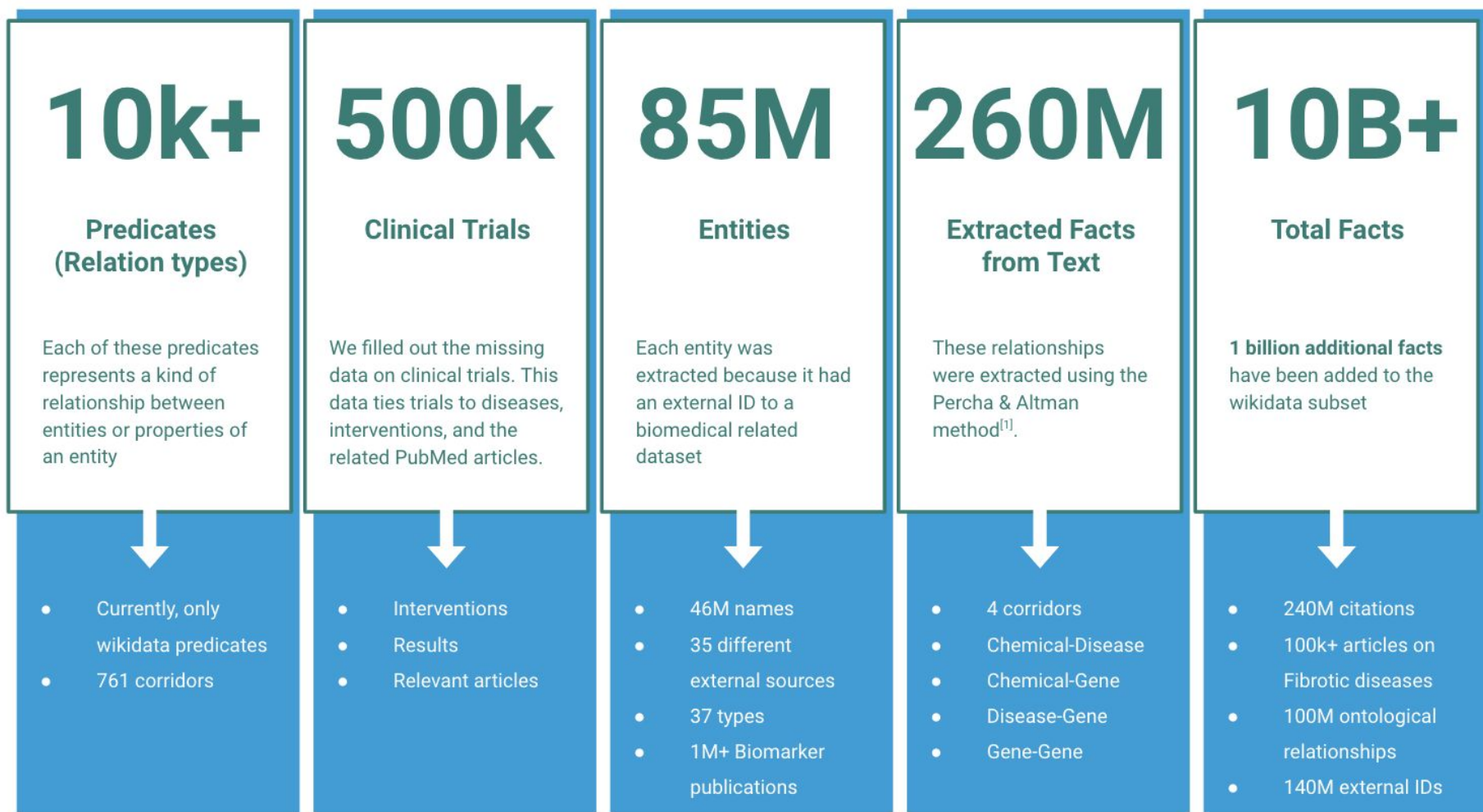


Figure 2. Summary of the knowledge graph contents

These counts are expected to increase as we add other datasets or increase what is extracted from current datasets.

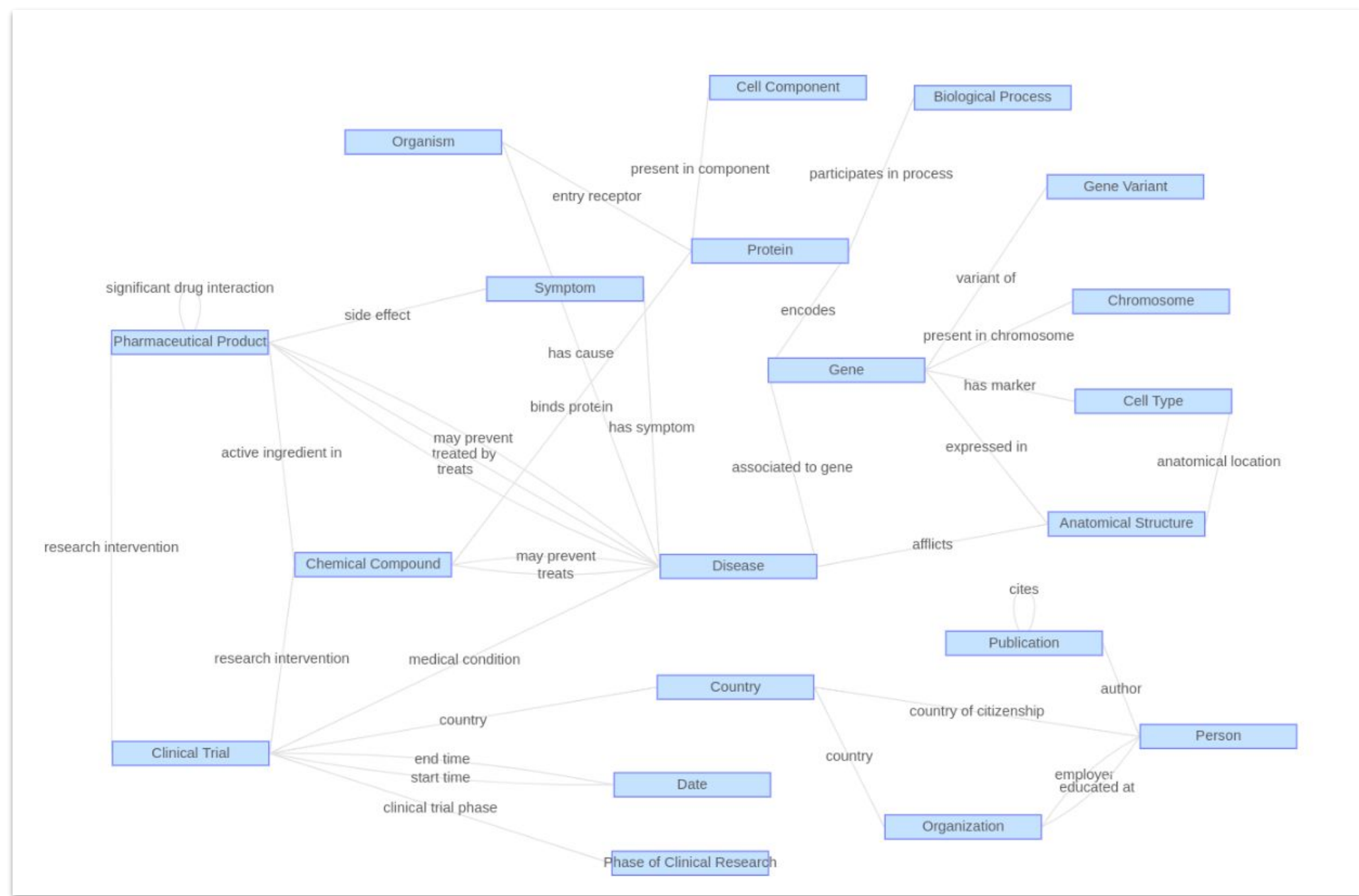
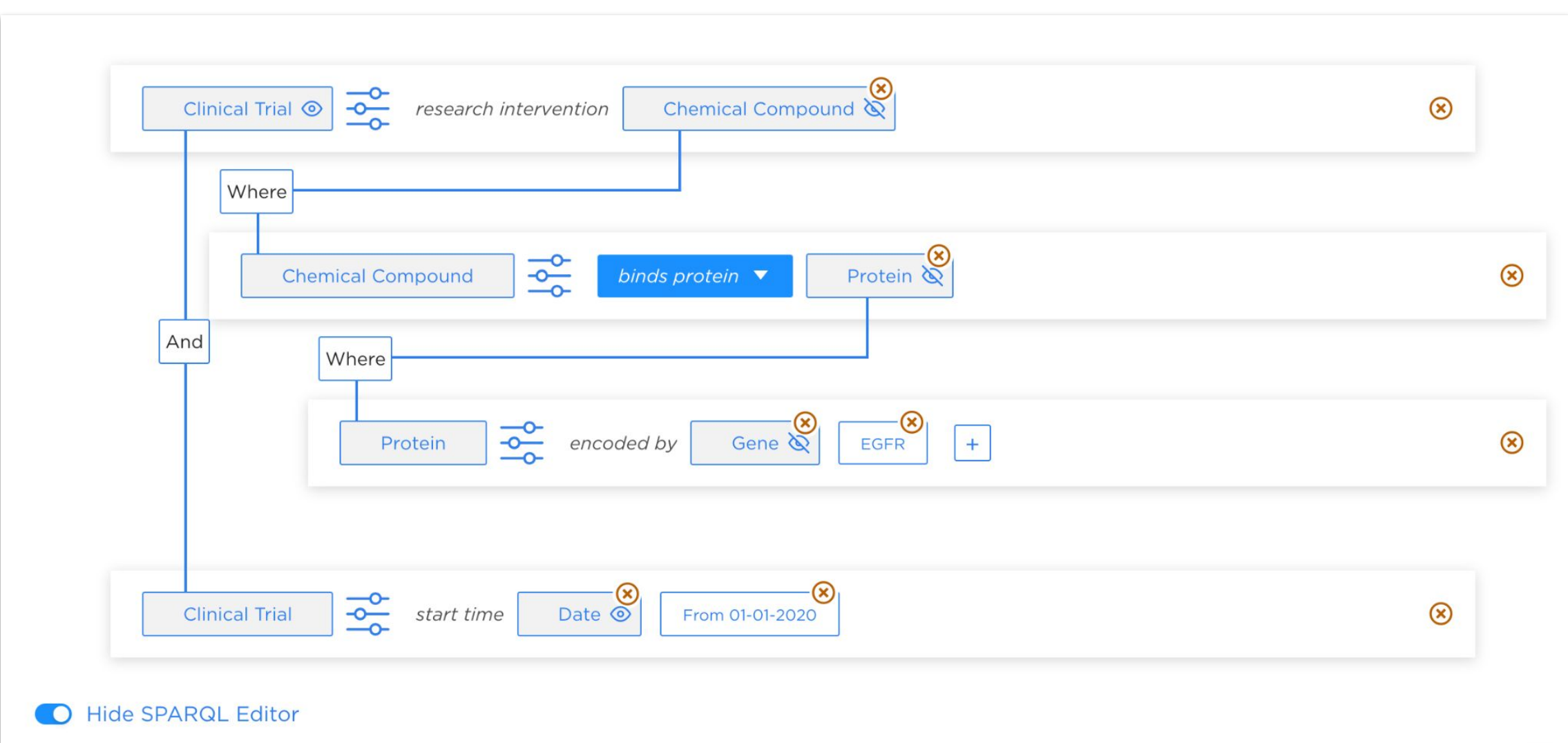


Figure 3. A sample of the schema of the knowledge graph

This visualization represents the entities and relationships of the graph subset currently accessible through the Visual Query Builder tool.

4.A. Query: Clinical trials since 01/2020 about compounds that bind to EGFR



4.B.

```
1 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 SELECT DISTINCT ?trial ?trial_label ?date_1 ?chemical_compound_1 ?chemical_compound_label ?protein_1 ?protein_label ?gene_2 ?gene_label WHERE {
4   ?trial <http://www.wikidata.org/prop/direct/P313> ?chemical_compound_1
5   ?trial <http://www.wikidata.org/prop/direct/P313> ?chemical_compound_label
6   ?chemical_compound_1 <http://www.wikidata.org/prop/direct/P313> ?protein_1
7   ?protein_1 <http://www.wikidata.org/prop/direct/P313> ?protein_label
8   ?protein_label <http://www.wikidata.org/prop/direct/P313> ?gene_2
9   ?gene_2 <http://www.wikidata.org/prop/direct/P313> ?gene_label
10  ?trial <http://www.wikidata.org/prop/direct/P313> ?date_1
11  ?date_1 <http://www.wikidata.org/prop/direct/P313> ?date_label
12  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
13  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
14  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
15  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
16  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
17  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
18  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
19  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
20  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
21  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
22  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
23  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
24  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
25  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
26  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
27  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
28  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
29  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
30  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
31  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
32  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
33  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
34  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
35  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
36  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
37  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
38  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
39  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
40  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
41  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
42  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
43  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
44  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
45  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
46  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
47  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
48  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
49  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
50  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
51  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
52  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
53  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
54  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
55  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
56  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
57  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
58  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
59  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
60  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
61  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
62  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
63  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
64  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
65  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
66  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
67  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
68  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
69  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
70  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
71  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
72  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
73  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
74  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
75  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
76  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
77  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
78  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
79  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
80  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
81  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
82  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
83  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
84  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
85  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
86  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
87  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
88  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
89  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
90  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
91  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
92  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
93  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
94  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
95  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
96  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
97  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
98  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
99  ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
100 ?date_label <http://www.wikidata.org/prop/direct/P313> ?date_label
```

Figure 4. Examples of querying the graph

4.A. This is an example of the Visual Query Builder. Using this, someone can easily query data in the graph through a series of drop-downs. This means that the user does not have to learn a new query language to access the knowledge graph.

4.B. This is the SPARQL (graph query language) generated by the Visual Query Builder. Users familiar with SPARQL already, they can write their own queries.

This	This Label	Date 1	Chemical Compound 1	Chemical Compound 1 Label	Protein 1	Protein 1 Label	Gene 2	Gene 2 Label
http://www.wikidata.org/	The Borneo Trial From History to Target: the Road to Personalized Target Therapy and Immunotherapy	2020-10-07T00:00:00.000Z	http://www.wikidata.org/	lapatinib	http://www.wikidata.org/	Epidermal growth factor receptor	http://www.wikidata.org/	EGFR
http://www.wikidata.org/	The Borneo Trial From History to Target: the Road to Personalized Target Therapy and Immunotherapy	2020-10-07T00:00:00.000Z	http://www.wikidata.org/	lapatinib	http://www.wikidata.org/	epidermal growth factor receptor	http://www.wikidata.org/	EGFR
http://www.wikidata.org/	LAT for Oligoprogressive NSCLC Treated With First-Line Osimertinib	2020-02-07T00:00:00.000Z	http://www.wikidata.org/	osimertinib	http://www.wikidata.org/	Epidermal growth factor receptor	http://www.wikidata.org/	EGFR
http://www.wikidata.org/	LAT for Oligoprogressive NSCLC Treated With First-Line Osimertinib	2020-02-07T00:00:00.000Z	http://www.wikidata.org/	osimertinib	http://www.wikidata.org/	epidermal growth factor receptor	http://www.wikidata.org/	EGFR

Figure 5. Tabular results of the query in Figure 4.

The query returns both the names, and entity URIs in a tabular format.

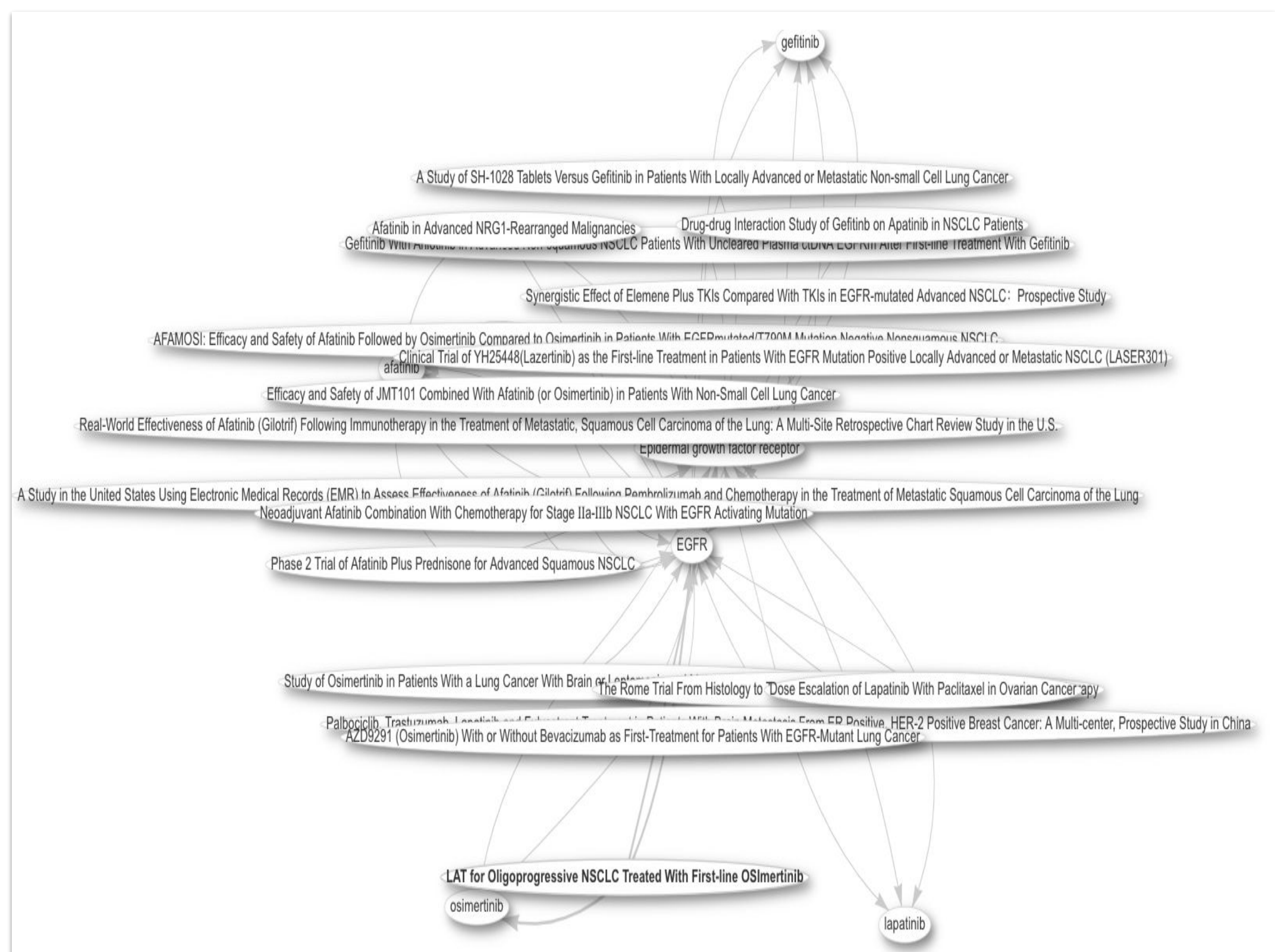


Figure 6. The graph visualization of the results in Figure 4.

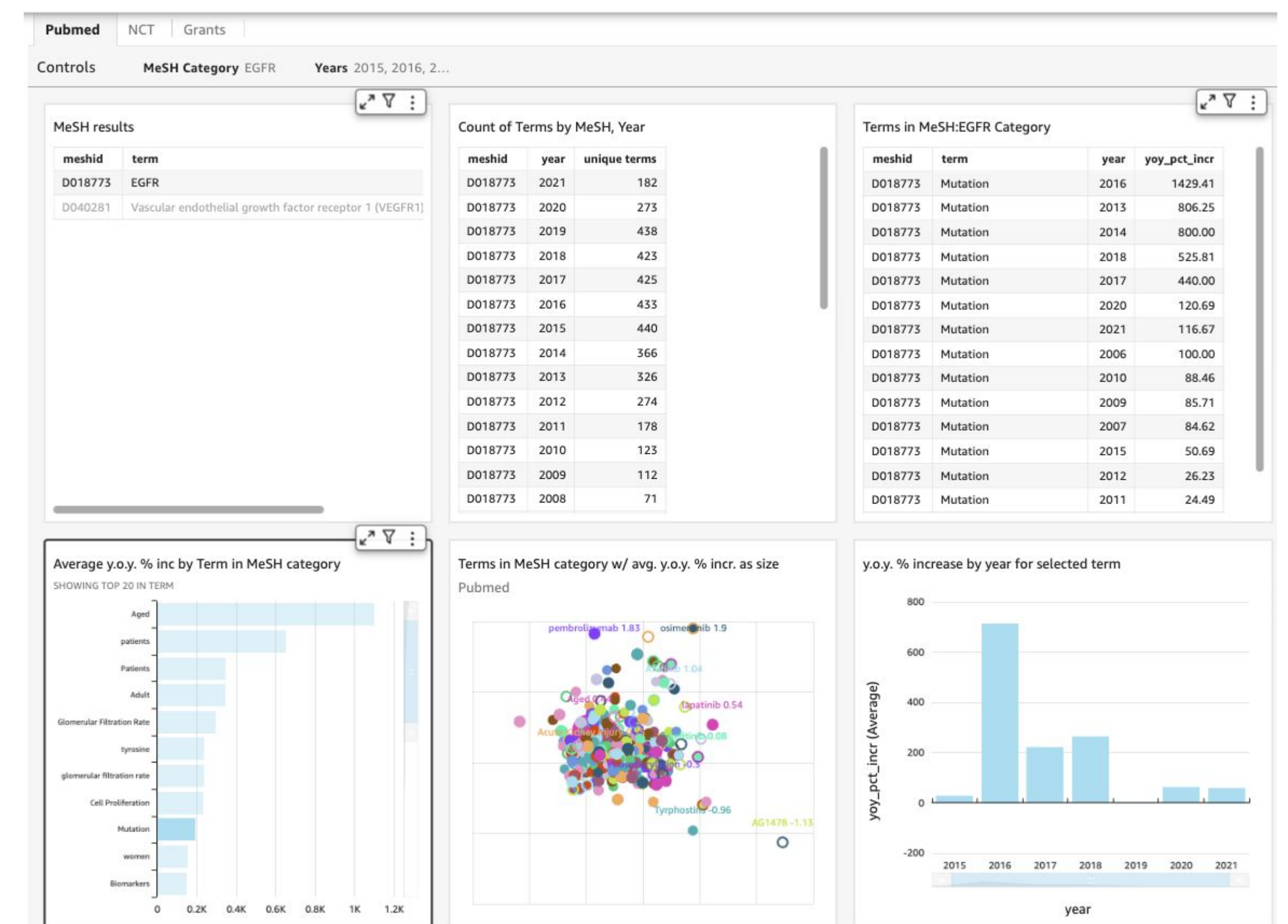


Figure 7. Trend Discovery analysis using MeSH terminology

Visualization of EGFR trend results in the Knowledge Graph across PubMed. "Mutation" term associated with EGFR and its change year over year for 2015-2020 is illustrated. Similar trends can be explored within Clinical Trials and NIH Grants information.

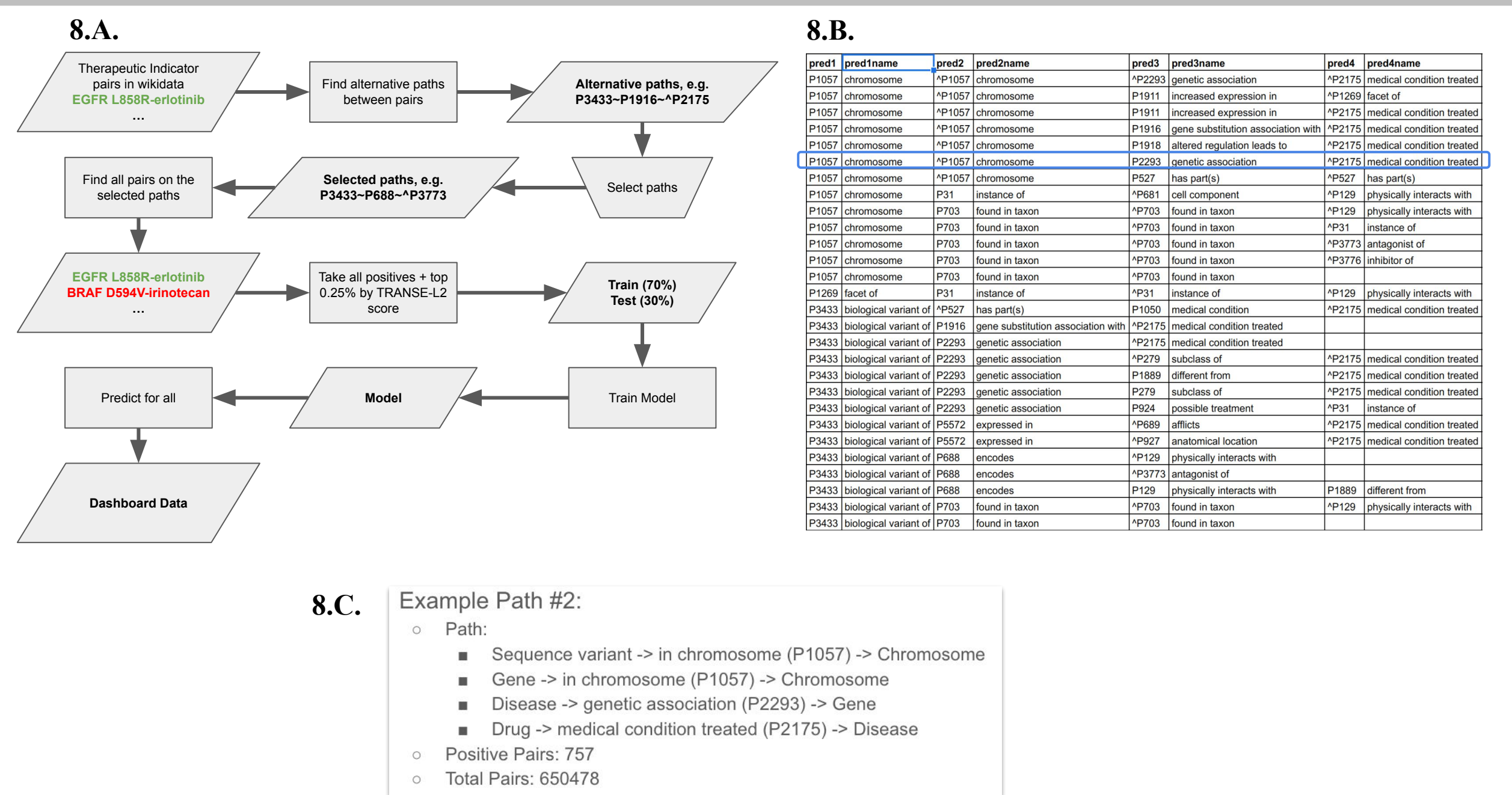


Figure 8. Gene Variant-to-Drug Link Prediction Task to Predict "Positive therapeutic Predictor" and "Negative Therapeutic Predictor" relationships in the Knowledge Graph: the alternative paths found, and used to generate the ranker training set 8.A. A table of all the paths found from sequence variant to drug. Note "P####" means inverted relationship 8.B. A high level description of the process to train the ranker model. 8.C. Example of a single path with the counts of positive gene variant-to-drug pairs, as well as total pairs identified

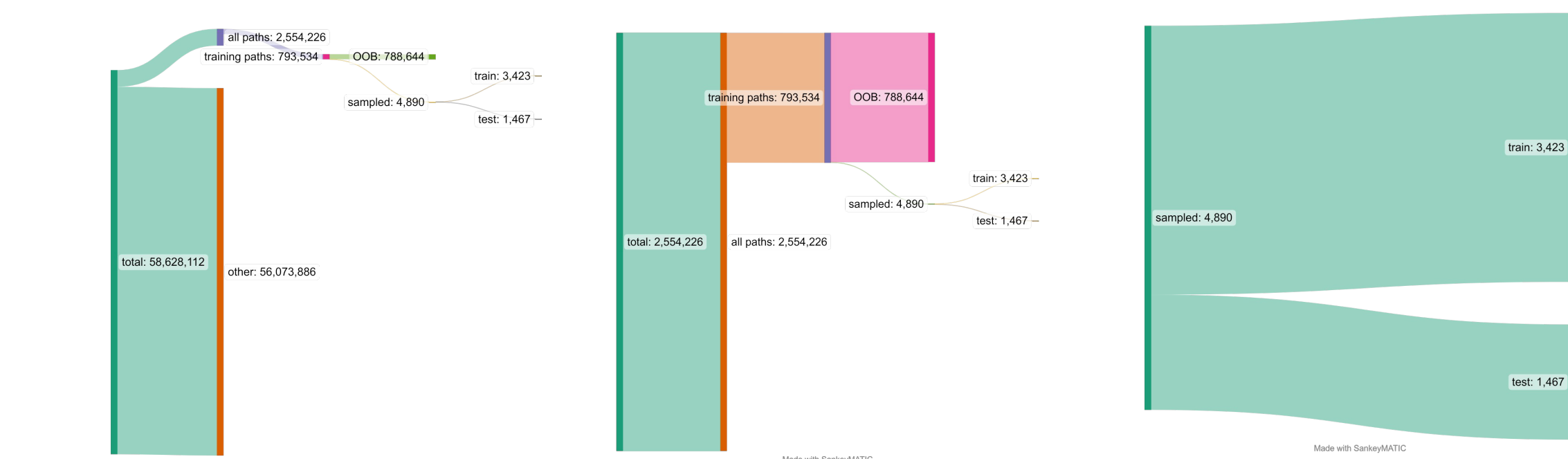


Figure 9. Constraining the data for training a ranker algorithm for predicting positive and negative therapeutic outcomes

There are 58 million possible sequence variant pairs. This is reduced by only selecting those pairs related by alternative paths, and then further sampling

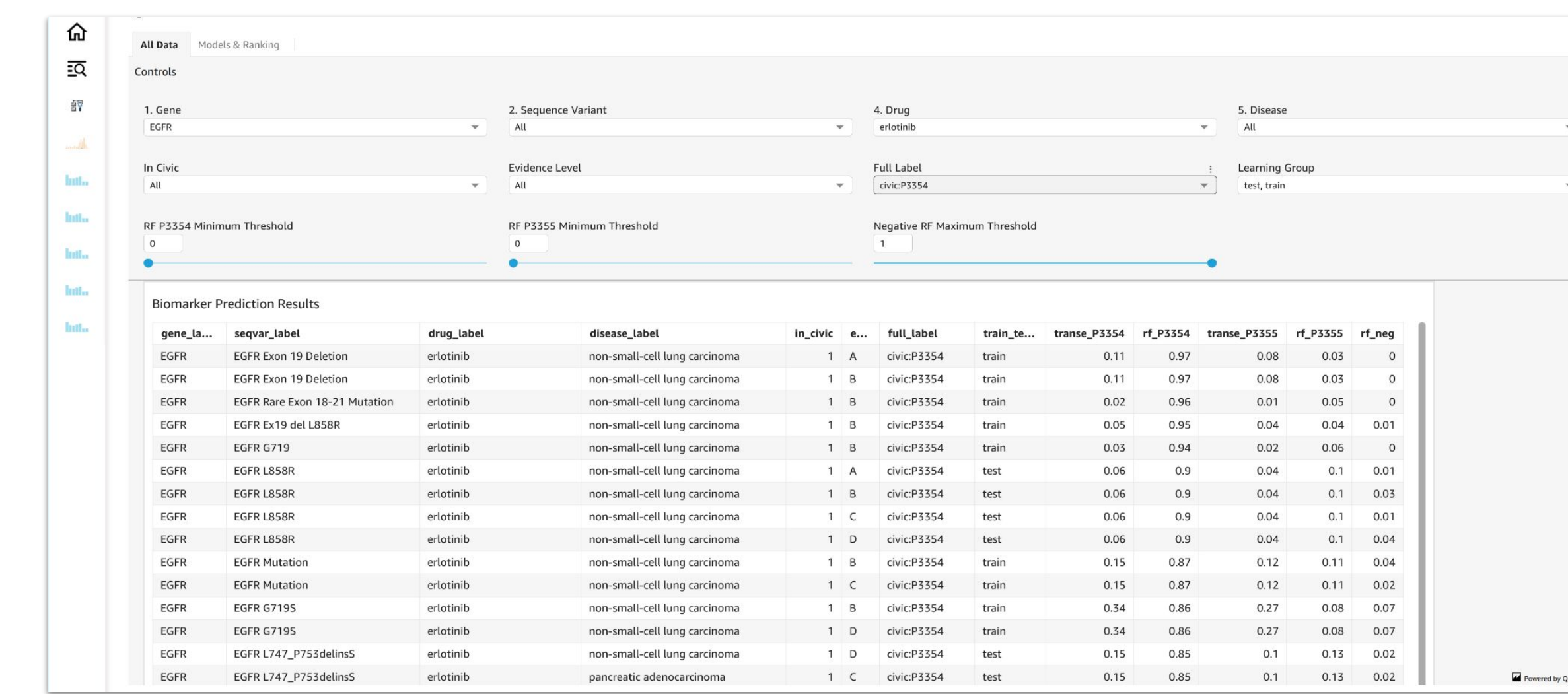


Figure 10. Examples of EGFR variants and Erlotinib drug relationship "P3354", denoting positive therapeutic prediction or predicting sensitivity/response to Erlotinib in non-small cell lung cancer.

Shown are examples or gene variants annotated in CIVIC database as predictive of sensitivity/response to Erlotinib in NSCLC that were reserved in the test set during algorithm training have been predicted correctly by the ranker, and score very similarly to the examples on which the training was done.

Conclusions

We demonstrate how combining AI innovations in NLP and graph analytics, as well as novel approaches in aggregating and harmonizing disparate sources of biomedical knowledge can act as a novel and promising digital solution with potential to accelerate biomedical knowledge, answer queries, discover important trends or assist in generating new ideas, and how knowledge graphs can be used for various medical purposes such as clinical decision support and drug discovery.

References

- Science Forum: Wikidata as a knowledge graph for the life sciences <https://doi.org/10.7554/elife.52614>
- Learning the Structure of Biomedical Relationships from Unstructured Text <https://pubmed.ncbi.nlm.nih.gov/26219079/>
- Knowledge Graph Embedding for Link Prediction: A Comparative Analysis <https://arxiv.org/abs/2002.00819>
- Finding melanoma drugs through a probabilistic knowledge graph <https://peerj.com/articles/cs-106/>
- Translating embeddings for modeling multi-relational data <https://dl.acm.org/doi/10.5555/2999792.2999923>
- <https://civicedb.org/evidence/home>

For additional information please contact: antoaneta.vladimirova@roche.com