



Large Language Model-Based Protein Sequence and Annotation Analysis System Developed for UniProt

Melike Akkaya, Rauf Yanmaz, Sezin Yavuz

Biological Data Science Laboratory, Department of Computer Engineering,
Hacettepe University, 06800, Ankara, Turkey

Advisor: Prof. Dr. Tunca Doğan

Introduction

The need to query large biological databases effectively and intuitively in biomedical research is steadily increasing. Although traditional query tools offer accuracy, they limit user accessibility due to the high level of technical expertise required.

"We have always observed that users, especially those new to UniProt, tend to explore the database using natural language queries... By understanding the meaning of these queries and 'translating' them into advanced queries, this capability will help users reach the expected results."

— Aurélien Luciani (EMBL-EBI)

In this project, a user-friendly system based on a large language model (LLM) has been developed, enabling natural language querying of the UniProt¹ protein sequence and annotation database.

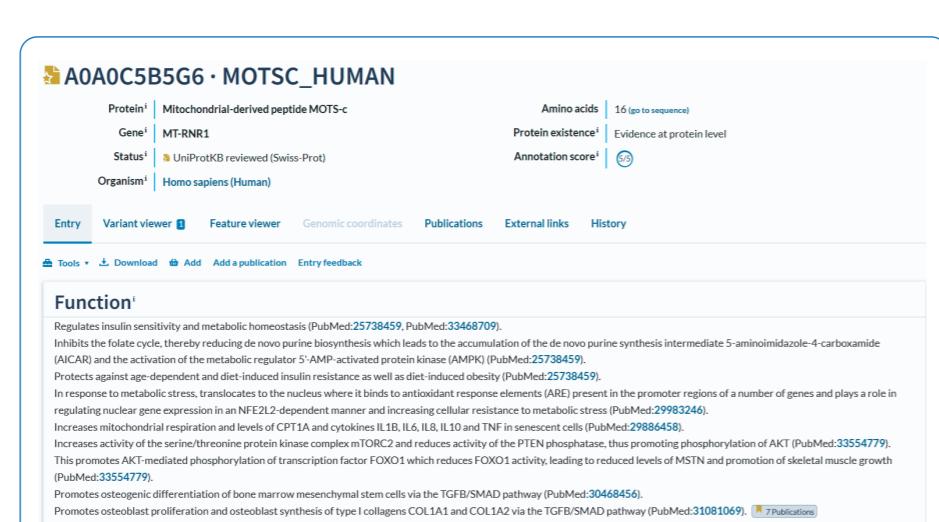


Figure 1: Example of Protein Information Page

Data

The system is composed of three main components, each supported by specific data types and sources:

1) LLM-Based Query Transformation:
User queries in natural language are adapted to UniProt's structure using documentation on query, search, and result fields as guidance.

2) Semantic Search:
Precomputed ProtT5 embeddings² from the UniProt Reviewed (Swiss-Prot) dataset (per_protein.h5) are used. Protein descriptions are taken from the uniprot_sprot.fasta file. For GO term analysis, only curated terms (excluding those labeled with IEA) from the go_basic.obo³ file are included.

3) Retrieval-Augmented Generation (RAG):
To generate context-aware explanations, the plain text uniprot_sprot.dat file is used. Since sequence embeddings are directly utilized, the raw sequence data is not separately processed.

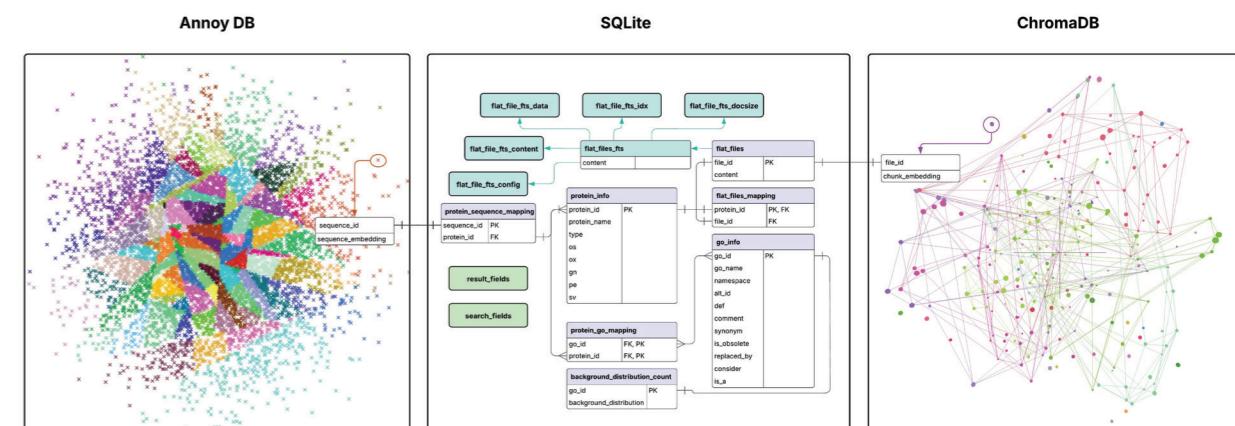


Figure 3: ER diagram of the databases used in the system

1) LLM-Based Solr Query Generation:

User queries in natural language are automatically converted into Solr queries for the UniProt database.

- Prompt Engineering:** UniProt's query, search, and result fields were analyzed to enable Solr query generation by LLMs. Search and result data were moved to an **SQLite-based helper database** to optimize prompts within token limits.

- Model Selection & Strategy:** Advanced LLMs were tested. Query variability was reduced using multiple trials and adjustable temperature settings, offering flexibility to users.

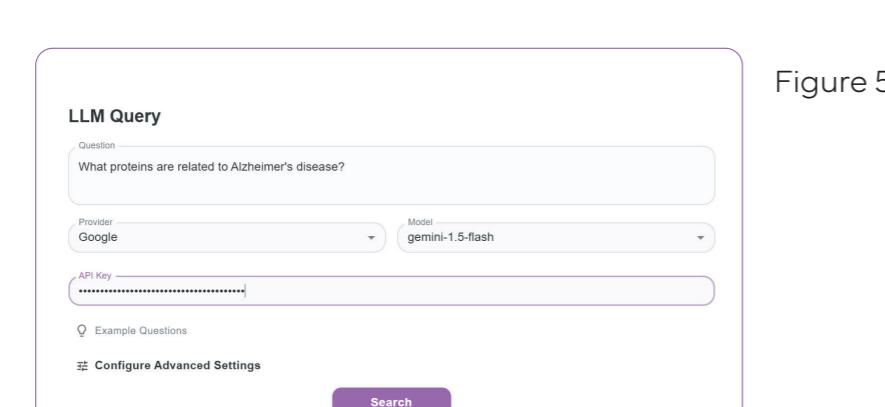
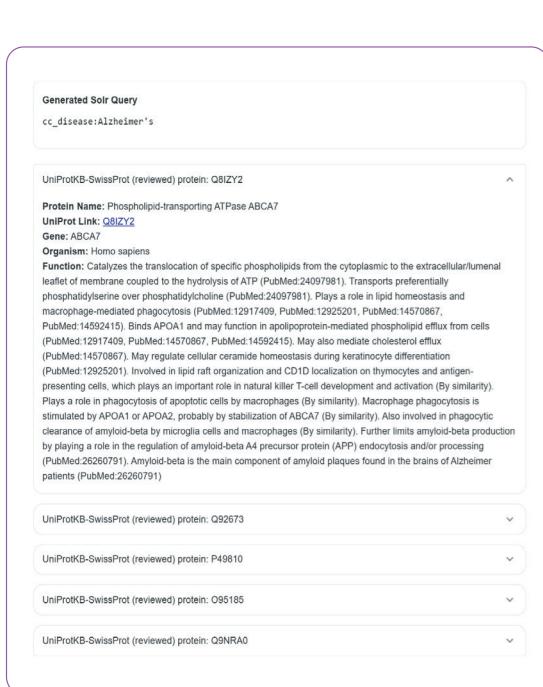
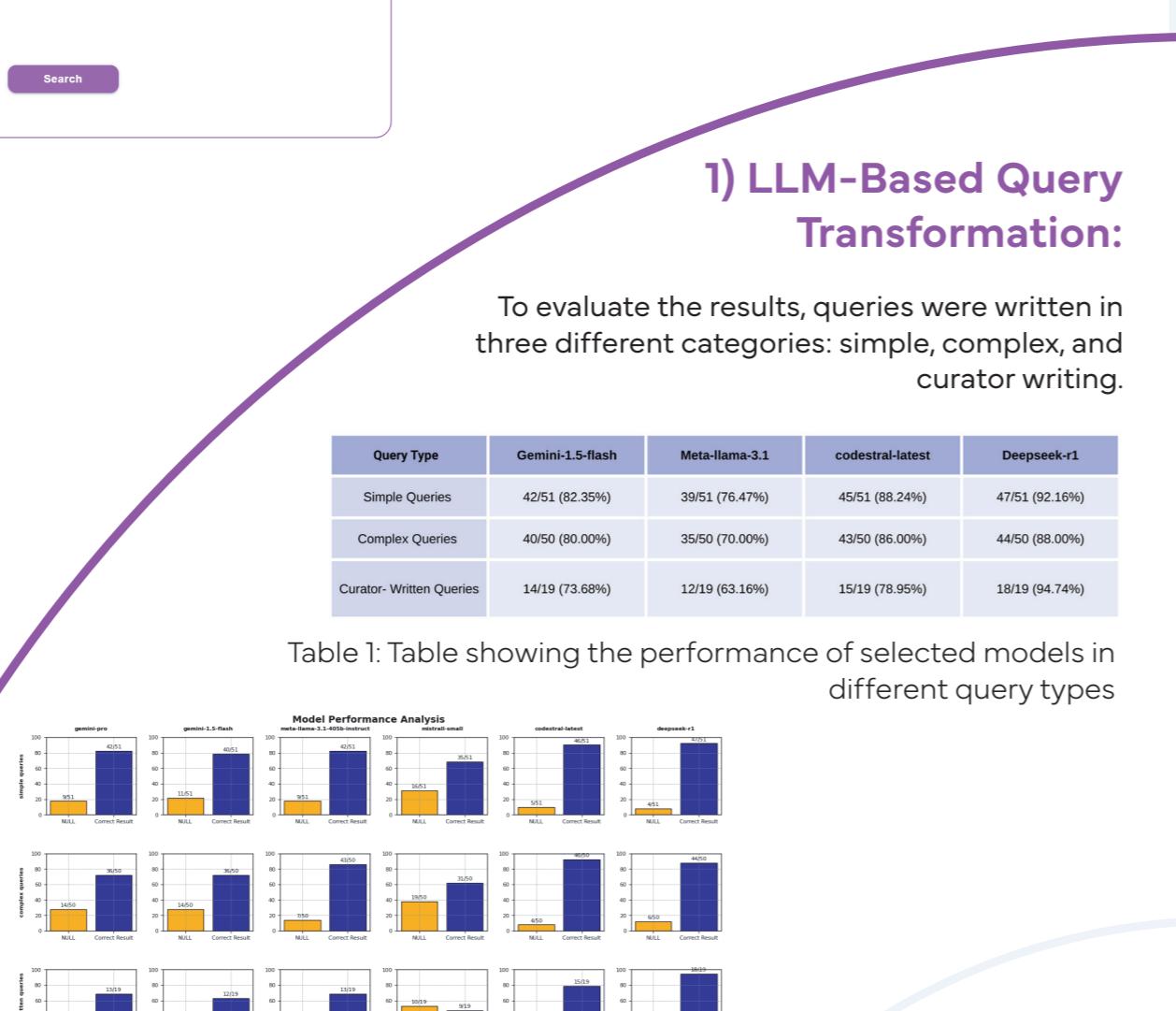


Figure 6

Users can interact with the system simply by writing questions in natural language, without needing to know the structure of the database or any query language (Figure 5).

The system translates these questions into technical queries suitable for UniProt, then analyzes the results and presents the most accurate answers in a clear and understandable way (Figure 6).



To evaluate the results, queries were written in three different categories: simple, complex, and curator writing.

Table 1: Table showing the performance of selected models in different query types

Query Type	Gemini-1.5-flash	Meta-Item-3.1	codestar-latest	Deepseek-v1
Simple Queries	42/51 (82.35%)	39/51 (76.47%)	45/51 (88.24%)	47/51 (92.16%)
Complex Queries	40/50 (80.00%)	35/50 (70.00%)	43/50 (86.00%)	44/50 (88.00%)
Curator-Written Queries	14/19 (73.68%)	12/19 (63.16%)	15/19 (78.95%)	18/19 (94.74%)

Figure 2 shows how difficult it is to understand the content of the data example used when setting up the vector database at this stage.

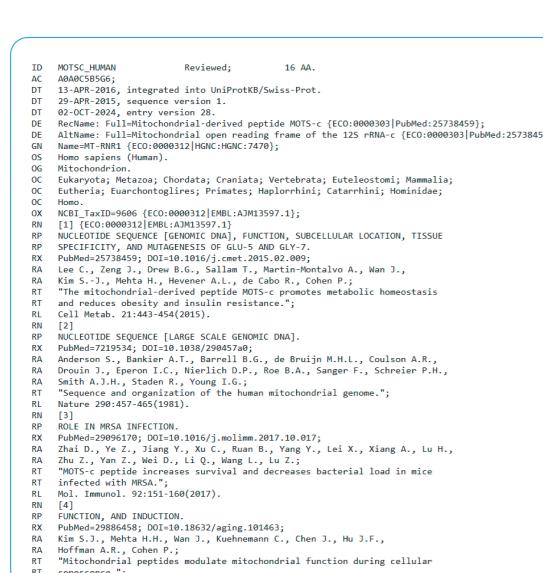


Figure 2: Sample flat file content

Figure 10 shows that users can optionally provide a protein sequence along with their question, without needing to know the query structure.

Figure 11 presents the resulting list of the most relevant and logical answers generated in response to such a query.

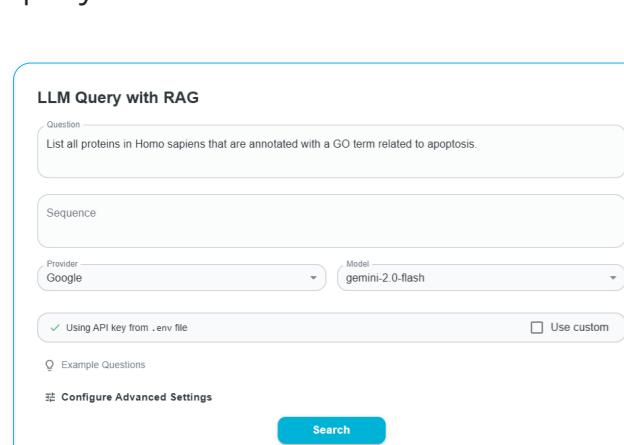


Figure 10

Retrieved Proteins						
Protein ID	Short Name	Protein Name	Organism	Taxon ID	Gene Name	pe SV
Q96188	HRK_HUMAN	Activator of apoptosis homolog	Homo sapiens	9606	HRK	1 1
Q96209	BIRC8_HUMAN	Bcl-2/retinoblastoma-binding protein 8	Homo sapiens	9606	BIRC8	1 2
Q96170	XXR9_HUMAN	XXL-related protein 9	Homo sapiens	9606	XXR9	1 1
Q96175	CSRNP2_HUMAN	Cysteine-rich secretory nuclear protein 2	Homo sapiens	9606	CSRNP2	1 1
Q96211	DODA_HUMAN	DNA damage-induced apoptosis suppressor protein	Homo sapiens	9606	DODA	1 2
Q96243	IWW12_HUMAN	Human-like 12	Homo sapiens	9606	MTRNR2L12	3
Q96281	IWW13_HUMAN	Human-like 13	Homo sapiens	9606	MTRNR2L13	3 1
Q96239	YGG41_HUMAN	Putative uncharacterized protein PNAS-138	Homo sapiens	9606	PNAS-138	5 1

Figure 11

3) RAG (Retrieval-Augmented Generation) Component:

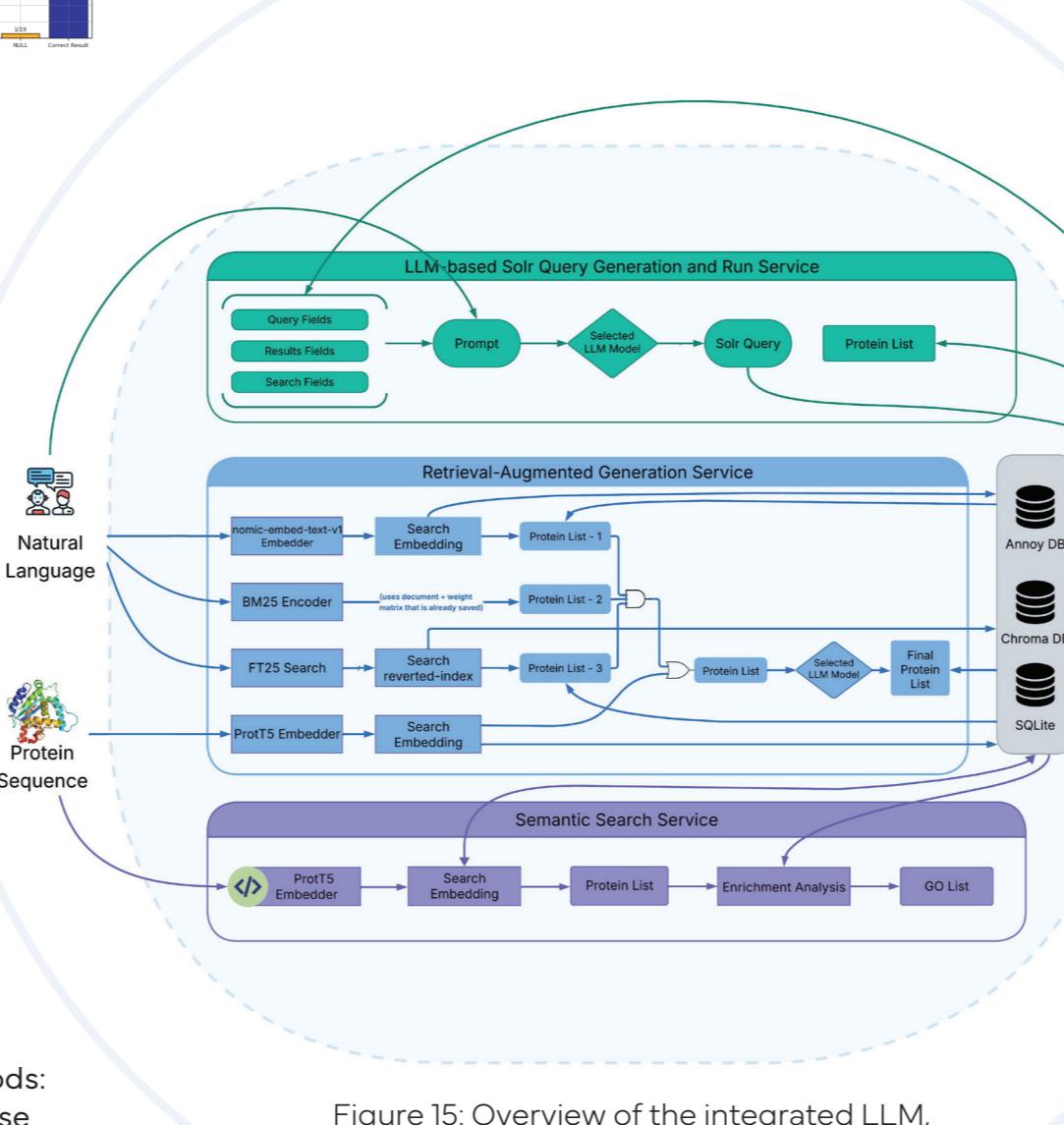


Figure 15: Overview of the integrated LLM, RAG, and semantic search services for

To compare a given protein sequence with others in the UniProt database in a biologically meaningful and fast way, a vector-based search system was developed. This system uses 1024-dimensional embeddings generated by UniProt's ProtT5 protein language model.

These vectors are indexed using the **Annoy library**⁵ for fast approximate nearest neighbor searches. Angular distance is used for similarity, and cosine similarity for ranking results. Since Annoy cannot store metadata, an **SQLite database**⁶ is used to link embedding IDs with protein information. This setup offers a portable, efficient, and file-based search infrastructure.

- Similarity Queries:** The user-provided protein sequence is embedded using ProtT5 and queried via Annoy to retrieve the most similar proteins in UniProt. The results are presented along with biological data such as **functional descriptions, GO terms, organism, and subcellular location**.

$$\text{Enrichment Score} = \frac{m}{\frac{M}{N}}$$

N: Total number of genes
M: Number of genes associated with the relevant GO term
n: User-provided gene list
m: Number of genes in the user list associated with the GO term

P-value: Measures the likelihood that the observed frequency of a GO term happened by chance. It's calculated using a hypergeometric distribution based on how many genes in a random list would be linked to the GO term.

$$P = \sum_{k=m}^n \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}}$$

2) Semantic Search and Embedding-Based Matching

As shown in Figure 12, tests conducted on 1,000 randomly selected sequences were compared with results obtained using the BLAST⁴ (Basic Local Alignment Search Tool) on the GenBank database. In **96.8% of the cases**, the correct match was found in the first position, and in nearly all cases, it was among the top 7 results.

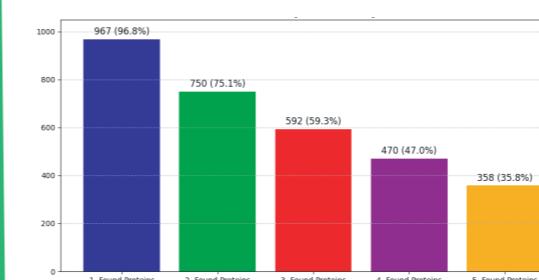


Figure 13: A chart showing the rank position of the correct protein match in the results for a sample query.

Two key metrics were used to evaluate system performance:

1) Spearman Correlation Coefficient:

Measures the similarity in ranking between vector-based search results and BLAST results. A value closer to 1 indicates stronger consistency between the two methods.

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

2) Common Protein Metric:

Measures the proportion of shared proteins identified by both methods. Ranges between 0 and 1; higher values indicate greater similarity.

$$\text{Common Protein Metric} = \frac{2 \times \text{Number of Common Proteins}}{\text{Proteins from Vector Search} + \text{Proteins from BLAST Search}}$$

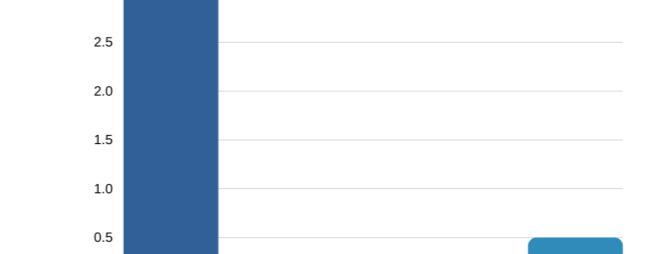


Figure 14: A graph showing the average speed performance of the system based on tests with 10,000 sequences, compared to currently used alternatives.

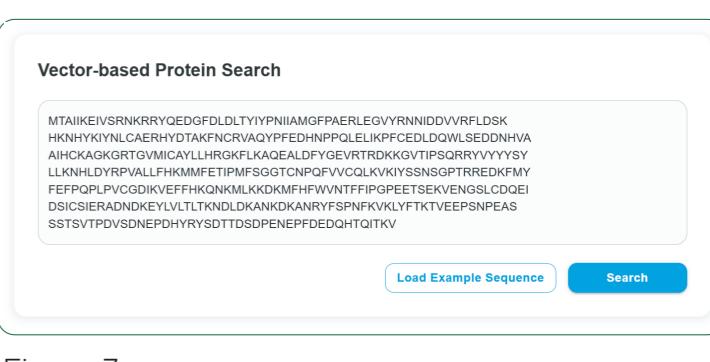


Figure 7



Figure 8

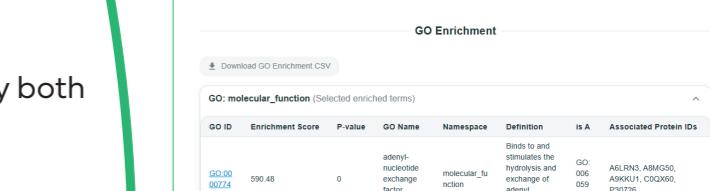


Figure 9

Figure 7 shows a system where the user can perform a search using only the available sequence information.

Figure 8 presents a list of proteins that are most similar to the given sequence and are therefore likely to be related.

Figure 9 displays the Gene Ontology (GO) terms associated with the relevant protein, grouped by their respective namespaces for easier interpretation by the user.

General Evaluation

This study presents a comprehensive system that can convert natural language queries in the UniProt database into Solr-based structural searches and perform semantic similarity analysis on protein sequences. This structure, which is faster and more user-oriented compared to traditional tools like BLAST, facilitates access to protein function information for both technical and non-technical users.

The system was developed in collaboration with the EMBL-EBI UniProt team, in line with real user needs and query habits.

Acknowledgements

We would like to thank Maria-Jesus Martin, the Protein Function Development team leader, team members Vishal Joshi and Aurelien Luciani, and the entire UniProt team and EMBL-EBI.

Important Links: <https://proquest.ngrok.app>

