

Cross Verification of Protein Information

Group 2

9th February 2021

Hassan Elsayed
Aya Abdelbaky
Daria Podorskaja
Mostafa Elhawwary

Outline

- Introduction (Hassan)
- Methodology (Aya)
- Software Overview (Daria)
- Discussion (Mostafa)

1. Introduction

Project overview

Cross Verification of Protein Information

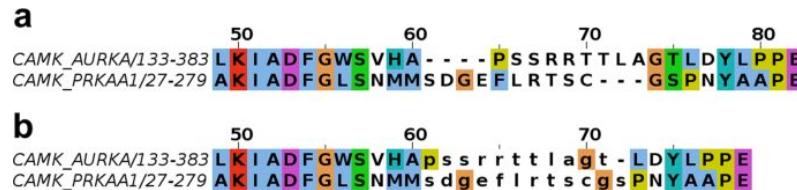
- Investigate whether protein information available on two different "gold standard databases" ; **Uniprot** and **NCBI refseq** match up
- Check whether the extracted comparable informations from the two databases are properly **updated**



RefSeq

Scientific significance

- Generate a **user friendly visual interface** of the data from Uniprot and Refseq
- To provide a **trustworthy package** for scientists & researchers to be able to retrieve non redundant and an up to date information about human proteins



Project problem/task

- Databases are not the same / uncomparable
 - query for the same protein in two different database will not get the same response
- Different research groups update or submit protein structure data independently
 - Difference can be in annotation, sequence, variant or PTM
- Compare data for all human genome across databases is a challenge

Background

- **Uniprot** is a protein sequence and annotation database., it provide a high-quality and freely accessible resource for most of the publicly available proteins of the world
- **Uniprot** is so-called gold standard“ as it provide a highly curated and annotated information about proteins using their two part system (Swiss-Prot and TrEMBL)
- **Refseq** “is an integrated, non-redundant, well-annotated set of reference sequences including genomic, transcript, and **protein**”.
- Refseq built by NCBI



NCBI
RefSeq

Uniprot cross reference

- Uniprot cross reference with many databases.
- Each database is described by its name and abbreviation and a link to its web server is provided, as well as literature references where available.

Cross-referencesⁱ

Sequence databases

<input checked="" type="radio"/> EMBL	AF036760 mRNA. Translation: AAC36493.1.
<input type="radio"/> GenBank	S82504 Genomic DNA. No translation available.
<input type="radio"/> DDBJ	S82502 Genomic DNA. No translation available.
	U60523 mRNA. Translation: AAB40387.1.
	S82500 Genomic DNA. Translation: AAB37501.1.
IPI	IPI00202716.
RefSeq	NP_036646.1. NM_012514.1.
UniGene	Rn.217584. Rn.48840.

3D structure databases

PDB	Entry 1L0B	Method X-ray	Resolution (Å) 2.30	Chain A	Positions 1589-1817	PDBsum [»]
PDBe	054952.					
RCSB PDB	054952. Positions 1-103, 1591-1801.					
PDBj	ModBase	Search...				

Protein-protein interaction databases

STRING	054952.
--------	---------

Genome annotation databases

Ensembl	ENSRNOT00000028109; ENSRNOP00000028109; ENSRNOG1
GenelD	497672.
KEGG	rno:497672.
UCSC	NM_012514. rat.

Organism-specific databases

Cross-references

Web resources

IL2Rbase

IL2RA mutation db

SeattleSNPs

Sequence databases

Select the link X01057 mRNA Translation: CAA25525.1

destinations: X03131 X03138 Genomic DNA Translation: CAA26906.1

EMBL K03122 mRNA Translation: AAB59535.1 Sequence problems.

GenBank M11066 M11065 Genomic DNA Translation: AAA67527.1

DDBJ AY563103 Genomic DNA Translation: AAS55572.1

AL157395 Genomic DNA No translation available.

AL137186 Genomic DNA No translation available.

CH471072 Genomic DNA Translation: EAW86414.1

M15864 Genomic DNA Translation: AAA59162.1

BN000945 Genomic DNA Translation: CAK26553.1

CCDS CCDS7076.1

PIR A44186, UHHU2

RefSeq NP_000408.1, NM_000417.2

NP_001295171.1, NM_001308242.1

NP_001295172.1, NM_001308243.1

interleukin-2 receptor subunit alpha isoform 1 precursor [Homo sapiens]

NCBI Reference Sequence: NP_000408.1

[Identical Proteins](#) [FASTA](#) [Graphics](#)

Go to:

LOCUS NP_000408 272 aa linear PRI 15-DEC-2020
DEFINITION interleukin-2 receptor subunit alpha isoform 1 precursor [Homo sapiens].

ACCESSION NP_000408

VERSION NP_000408.1

DBSOURCE REFSEQ: accession NM_000417.3

KEYWORDS RefSeq; MANE Select.

SOURCE Homo sapiens (human)

ORGANISM Homo sapiens

Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominoidea; Homo.

REFERENCE 1 (residues 1 to 272)

AUTHORS Ni M, Tian FB, Xiang DD and Yu B.

TITLE Characteristics of inflammatory factors and lymphocyte subsets in patients with severe COVID-19

J Med Virol 92 (11), 2600-2606 (2020)

[32470153](#)

REFERENCE 2 (residues 1 to 272)

AUTHORS Stefanovic M, Zivotic I, Stojkovic L, Dincic E, Stankovic A and Zivkovic M.

TITLE The association of genetic variants IL2RA rs2104286, IFI30 rs11554159 and IKZF3 rs12946510 with multiple sclerosis onset and

3D structure databases

Select the link PDB entry

1ILM model

Method

-

Resolution (Å)

-

JOURNAL

J Med Virol 92 (11), 2600-2606 (2020)

PUBMED

[32470153](#)

REFERENCE

2 (residues 1 to 272)

AUTHORS

Stefanovic M, Zivotic I, Stojkovic L, Dincic E, Stankovic A and Zivkovic M.

TITLE

The association of genetic variants IL2RA rs2104286, IFI30

rs11554159 and IKZF3 rs12946510 with multiple sclerosis onset and

destinations:

PDBe

RCSB PDB

PDB

Program implementation/Solve tasks

- ★ Uniprot and refseq are cross reference to each others
- ★ Download all human proteins sequence from Uniprot and refseq
 - Send API query to extract data needed from databases
 - Ensure gene symbol, accession identifier, amino acid sequence are found in the downloaded data
- ★ Map human proteins from the two databases (choose a column)
- ★ Compare between uniprot and refseq data (comparable value)
 - accession IDs , gene symbol, amino acid sequence
- ★ Show the data in a web interface displaying a comparison between the two databases.

2. Methodology

Strategy

- **Task 1**
 - Database Download (Daria)
- **Task 2**
 - Mapping the Metadata (Hassan)
- **Task 3**
 - Compare the Data (Aya)
- **Task 4**
 - GUI (Mostafa)

Implementation

- Step 1
 - Download the entire Human Proteome from UniProt and NCBI RefSeq databases
 - Write the downloaded metadata to separate “.tsv” files
- Step 2
 - Merge both files using “Accession Number Version”
 - Write the merged metadata to a new “.tsv” file
- Step 3
 - Generate a dataframe with the specified protein metadata obtained from both databases
 - Generate a dataframe with the percentage similarity for across the whole Human Proteome
- Step 4
 - View Table of Summary Statistics for comparing the databases
 - View Table for the metadata of a protein given by the user

Competitors

Software Tool for
Researching Annotations of
Proteins (STRAP)



LAB / MaxQuant



customProDB /
R Package



Protein Metrics /
Featured Byos®
Workflows



Olson Lab /
Optide-Hunter

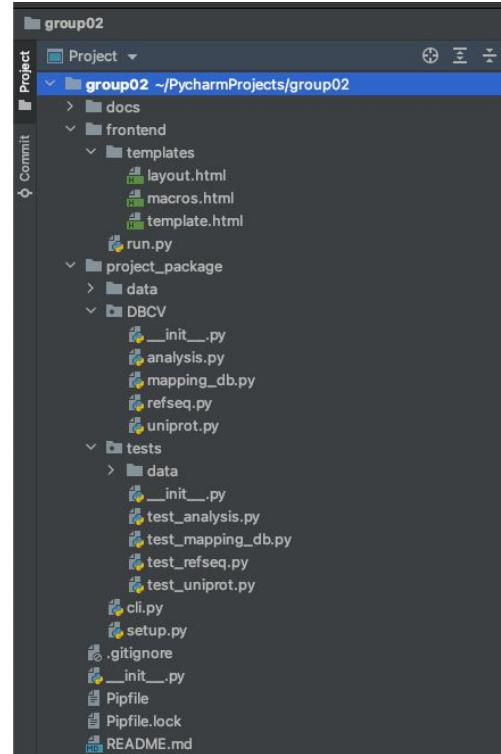


DBCV

3. Software Overview

Project structure

- Frontend
 - Templates
 - run.py
- Project package
 - Data
 - DBCV
 - Tests
 - Test data
- cli.py
- setup.py
- pipfile
- README.md



Download Data from RefSeq (refseq.py)

```
12     def download_refseq(output_path: str, entries_number: int = 120_000):
13         """
14             Fetch entries form NCBI protein database using Entrez API functions.
15             API query: "Homo sapiens[Orgn] AND refseq[filter]"
16             GenBank parse is used for reading the fetched data.
17             https://biopython.readthedocs.io/en/latest/chapter\_entrez.html
18             :param entries_number: number of entries
19             :param output_path: path to output file
20             :return:
21             """
22
23         if os.path.exists(output_path):
24             print(f"File {output_path} already exists")
25             return
26
27         add_headers(output_path=output_path)
28         handle = Entrez.esearch(db="protein", term="Homo sapiens[Orgn] AND refseq[filter]", retmax=entries_number)
29         record = Entrez.read(handle)
30         for i in range(0, entries_number, 101):
31
32             handle = Entrez.efetch(db="protein", id=record['IdList'][i:i + 100], rettype="gp", retmode="text")
33             records = GenBank.parse(handle)
34
35             for r in records:
36                 data = augment_data(r)
37                 write_data(data, output_path)
38         handle.close()
```

Download Data from Uniprot (uniprot.py)

```
10     def generate_payload():
11         """
12             Payload for the API query "https://www.uniprot.org/uniprot/?query=proteome:UP000005640"
13             Documentation can be found here:
14             https://bioservices.readthedocs.io/en/master/\_modules/bioservices/uniprot.html
15             :return: payload for the API query
16         """
17
18         return {
19             'query': 'organism :"Homo sapiens (Human) [9606]" AND proteome:up000005640',
20             'format': 'tab',
21             'columns': ','.join(['genes', 'id', 'database(HGNC)', 'database(RefSeq)', 'sequence']),
22         }
23
24     def download_uniprot(payload: dict, output_path: str = "../data/uniprot.tsv"):
25         """
26             Download information from UniProt based on the payload.
27             :param payload: payload for API query
28             :param output_path: path for the output file
29             :return:
30         """
31
32         if os.path.exists(output_path):
33             print(f"File {output_path} already exists")
34             return
35         result = requests.get(BASE + KB_ENDPOINT, params=payload)
36         if result.ok:
37             with open(output_path, 'a') as f:
38                 f.write(result.text)
39         else:
40             print(result.status_code)
```

Mapping of the Data (mapping_db.py)

```
60     # Merging the modified Uniprot dataframe and the RefSeq dataframes based on the RefSeq accession number
61
62     df4_mergefiles = pd.merge(refseq_df, df3_newuniprot, on='accession_number_version')
63
64     # Rearranging the columns of the dataframe for better interpretation
65
66     df4_mergefiles = df4_mergefiles[
67         ['accession_number_version', 'Refseq_accession_number_Refseq', 'db_source_accession_Refseq', 'Entry_Uniprot',
68          'Entry_Uniprot_Refseq',
69          'HGNC_Uniprot', 'HGNC_Refseq', 'MimID_Refseq', 'Gene_Names_Uniprot', 'Gene_Name_Refseq',
70          'Gene_synonyms_Refseq',
71          'Gene_ID_Refseq', 'Sequence_Uniprot', 'Sequence_Refseq']]
72
73     # Writing the dataframe to an output file
74     df4_mergefiles.to_csv(out_path, sep='\t', index=False, header=True)
```

Analysis of the Data (analysis.py)

```
37     protein_info_data = [
38         "Protein's gene symbol", str(df.iloc[row, df.columns.get_loc('Gene_Names_Uniprot')]).split(' ')[0],
39         df.iloc[row, df.columns.get_loc('Gene_Name_Refseq')],
40         str(df.iloc[row, df.columns.get_loc('Gene_Names_Uniprot')]).split(' ')[0] == df.iloc[
41             row, df.columns.get_loc('Gene_Name_Refseq')]],
42         ["UniProt accession IDs", df.iloc[row, df.columns.get_loc('Entry_Uniprot')]],
43         df.iloc[row, df.columns.get_loc('Entry_Uniprot_Refseq')],
44         df.iloc[row, df.columns.get_loc('Entry_Uniprot')] == df.iloc[
45             row, df.columns.get_loc('Entry_Uniprot_Refseq')]],
46         ["RefSeq accession IDs", df.iloc[row, df.columns.get_loc('accession_number_version')]],
47         df.iloc[row, df.columns.get_loc('Refseq_accession_number_Refseq')],
48         df.iloc[row, df.columns.get_loc('accession_number_version')].split('.')[0] == df.iloc[
49             row, df.columns.get_loc('Refseq_accession_number_Refseq')]],
50         ["Amino Acid sequences lengths", len(df.iloc[row, df.columns.get_loc('Sequence_Uniprot')]),
51         len(df.iloc[row, df.columns.get_loc('Sequence_Refseq')]),
52         len(df.iloc[row, df.columns.get_loc('Sequence_Uniprot')]) == len(
53             df.iloc[row, df.columns.get_loc('Sequence_Refseq')])],
54         ["Amino Acid sequences alignments", uniprot_seq,
55          refseq_seq,
56          df.iloc[row, df.columns.get_loc('Sequence_Uniprot')] == df.iloc[
57              row, df.columns.get_loc('Sequence_Refseq')]]
58     ]
59     protein_info = pd.DataFrame(data=protein_info_data, columns=protein_info_columns)
```

Frontend (run.py)

```
13     @app.route("/")
14     def home():
15         """
16             This method loads the home page with the comparison table between the two databases
17
18         """
19         df = compare_data('new_out.tsv')
20         df = df.to_html(index=False)
21         return render_template('template.html', DF=df)
22
23
24     @app.route("/protein", methods=['POST'])
25     def protein():
26         """
27             This method is responsible for displaying the comparison between the two databases for a specific protein written
28             in the textbox of the GUI
29
30         """
31         summary = compare_data('new_out.tsv')
32         summary = summary.to_html(index=False)
33         protein_name = request.form['textbox']
34         df = compare_data('new_out.tsv', protein_name)
35         if df is not None:
36             df = df.to_html(index=False)
37
38         return render_template('template.html', comp=df, DF=summary)
```

Encountered Problems

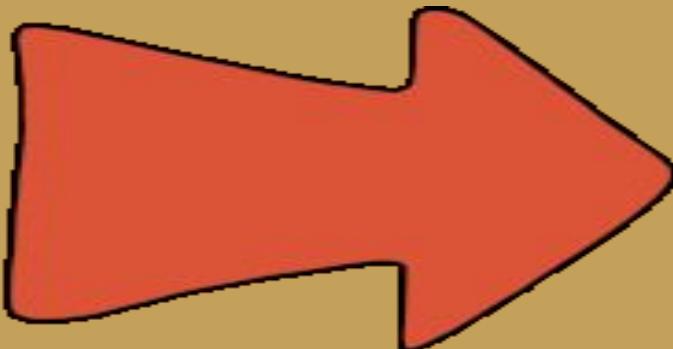
- Finding the right NCBI database
- Size of the downloaded files from the database
- Sequence is too long to display on frontend for comparison

4. Discussion

Use case:

- Researchers working on protein data
- Used to decide whether the info they have is consistent or not
 - Source?
 - Matching or not!?

Live Demo.



References

Google search. (n.d.). Retrieved February 20, 2019, from https://www.google.com/search?q=c%23%2B&tbm=isch&ved=2ahUKEwjMw6TCp9buAhUOy6QKHTWAnQ2-cCegQIABAA&oq=c%23%2B&gs_lcp=CgNpbWcQAzIECMQJzIECCAAyAggAMgQIABBDMgIADICCAyAggAMgIADICCABQ_cwgWP3AWCMwyFoAHAAeECAAAT-IAT-SAQExmAEAOAEBqgELZ3dzLdpel1pbVfAAQE&scilnt=1&imgt=1&imgi=1536&rlz=1C1SQJL_enDE872DE872#imgrc=f-B4RoxPF0wVoM&imgdii=54dhvB9W_Oxs

Google search, (n.d.). Retrieved February 06, 2021, from
https://www.google.com/search?q=optide-hunter&rlz=1C1GCEU_enUS848US848&tbo=q&sourceid=chrome&ie=UTF-8

Google search. (n.d.). Retrieved February 06, 2021, from
https://www.google.com/search?q=2Bpackage&tbm=isch&ved=2ahUKEwjo1Y7spbdBuAhUhrQKQH1dAQS2C-qegQIABAA&oq=r%2Bpackage&gs_lcp=CgNpbWcQAzIeCMQjzCCAAyAggAMglIADICCAyAggAMglIADICCAyAggAUN_QAVfOAEwMHP-AUhA-ALAEfGFLPRA7pA-aACpC2d2p12-YetA-M1wAEF8_0mavclipnTjim+0mavclipnVWV1Z7AVoVUH7EV_0mavclipn-7E42-mavclipn-15268_mavclipn-1C1S0L_oD5F72D5F72#immrc=01k1132pui1&opq=india&igc=1USldShBM