

# **Algorithms in Structural Bioinformatics**

## **Assignment #2**

**Konstantinos Giatras**

## Exercise 1

The 3D structure of the human serotonin receptor was determined in the absence and in the presence of serotonin (PDB id: 7e2x, 7e2y). For this exercise, we will be going through the information/figures provided in PDBe regarding the 7e2y entry (<https://www.ebi.ac.uk/pdbe/entry/pdb/7e2y>).

### Question 1.1

To find the **method** by which the 3D structure of 7e2y was determined, we download the pdb file and locate the appropriate information in the file. In this case, the structure was determined using **cryo-EM** experiments. This is further confirmed by the relevant paper (<https://www.nature.com/articles/s41586-021-03376-8>).

The **resolution** of the structure can be found either on the PDBe entry or the downloaded pdb file. For 7e2y, the resolution is **3.00 Angstroms**.

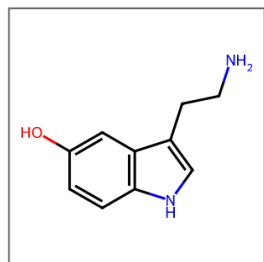
### Question 1.2

The **number of chains and amino acids in each chain** can be found under the structure analysis section of the 7e2y PDBe entry. They are as follows:

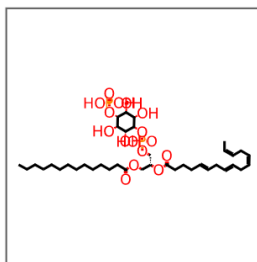
No	Chain	Amino Acids
1	A	354
2	B	345
3	G	71
4	R	543

### Question 1.3

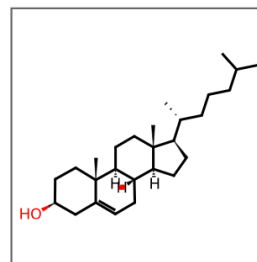
The **ligands** present in the structure can be found under the ligands and environments section of the 7e2y PDBe entry. These include 1x **serotonin** (SRO), 4x **cholesterols** (CLR) and 1x [(2R)-1-[oxidanyl-[(2R,3R,5S,6R)-2,3,5,6-tetrakis(oxidanyl)-4-phosphonooxy-cyclohexyl]oxy-phosphoryl]oxy-3-tetradecanoyloxy-propan-2-yl] (5E,8E)-hexadeca-5,8,11,14-tetraenoate (**J40**) molecule.



1 x SRO



1 x J40



4 x CLR

## Exercise 2

For this exercise, we will be using Chimera-X (<https://www.cgl.ucsf.edu/chimerax/>), a molecular visualization program, to prepare different figures of the overall 7e2y structure of the complex (this website: <https://www.cgl.ucsf.edu/chimerax/docs/user/commands/colortables.html>, was used to learn the appropriate coloring commands).

### Question 2.1



Figure 1: Secondary Structure Elements of the Protein Coloured by Chain (A: salmon, B: sky blue, G: purple, R: gold)

## Question 2.2

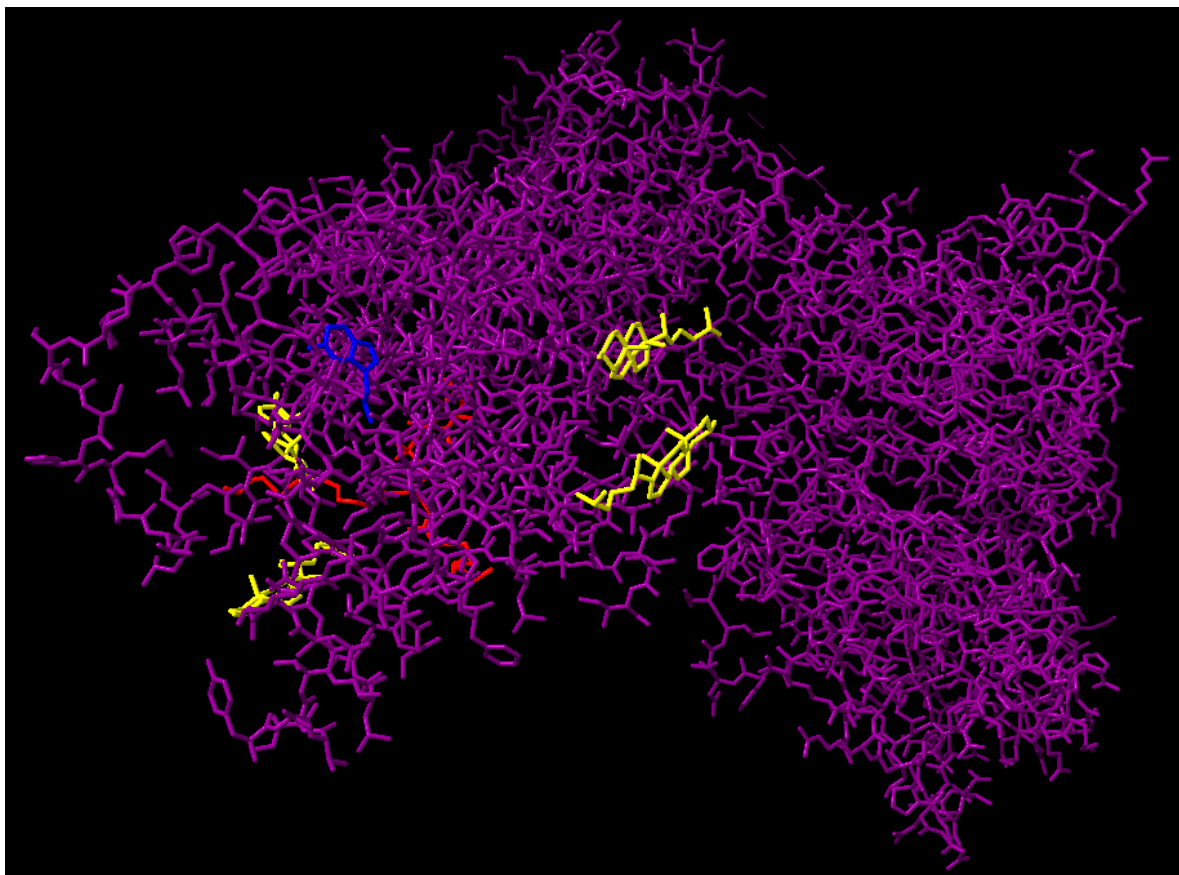


Figure 2: Chemically Distinct Molecules in the Structure (Human Serotonin Receptor: purple, Serotonin: blue, Cholesterol: yellow, [(2R)-1-[oxidanyl-[(2R,3R,5S,6R)-2,3,5,6-tetrakis(oxidanyl)-4-phosphonooxy-cyclohexyl]oxy-phosphoryl]oxy-3-tetradecanoyloxy-propan-2-yl] (5E,8E)-hexadeca-5,8,11,14-tetraenoate: red)

## Exercise 3

For this exercise (based on this article: <https://www.ebi.ac.uk/pdbe/news/serotonin-ruffling-feathers>), we will determine the c-RMSD between the human (7e2y, deposited in PDB, determined experimentally) and chicken (D2K8P9, predicted using AlphaFold, <https://alphafold.ebi.ac.uk/entry/D2K8P9>) serotonin 1A receptor structures (5-HT1A), mainly using commands on the console of the PyMOL software, which provides a Python API for structural analysis.

### Question 3.1

To determine the RMSD over all atoms, we begin by loading the two pdb files into PyMOL. We then type the following commands:

```
1 align 7e2y, AF-D2K8P9-F1-model_v4, cycles=0, object=aln
2 rms_cur 7e2y & aln, AF-D2K8P9-F1-mod & aln, matchmaker=-1
```

The *'align'* command (<https://pymolwiki.org/index.php/Align>) performs a sequence alignment followed by a structural superposition, and then carries out zero or more cycles of refinement in order to reject structural outliers found during the fit. The all-atom RMSD can be obtained by setting *'cycles=0'* and thus not doing any outlier rejection. The *'object'* argument creates and names an alignment object to be used by the next command.

The *'rms\_cur'* command ([https://pymolwiki.org/index.php/Rms\\_cur](https://pymolwiki.org/index.php/Rms_cur)) computes the RMS difference between two atom selections without performing any fitting. By default, only matching atoms in both selections are used for the fit. The previously generated alignment object is used in the atom selection and the *'matchmaker=-1'* argument (for matching atom pairs) is used to assume that atoms are stored in the identical order.

The output in the console is: **all-atom RMSD = 9.131 Angstrom (2162 to 2162 atoms)**, indicating a relatively high level of structural deviation between the two proteins, with significant differences in the positions of all atoms. The comparison is subject to possible errors, since it was conducted between predicted and experimentally determined structures.

### Question 3.2

To determine the RMSD over  $C_{\alpha}$  atoms, we begin by loading the two pdb files into PyMOL. We then type the following commands:

```
1 align 7e2y, AF-D2K8P9-F1-model_v4, object=aln
2 rms_cur 7e2y & aln and name ca, AF-D2K8P9-F1-mod & aln and name ca, matchmaker=-1
```

In the *'align'* command, we remove the *'cycles'* argument and in the *'rms\_cur'* command, we add two *"and name ca"* arguments to calculate the RMSD only for the  $C_{\alpha}$  atoms.

The output in the console is:  **$C_{\alpha}$  RMSD = 1.100 Angstrom (235 to 235 atoms)**, indicating a relatively higher degree of similarity in the backbone conformation of the two proteins. However, there may still be variations in the side chain orientations and overall protein structure.

### Question 3.3

To determine the RMSD using the secondary structure elements only, we need to define the secondary structure elements (SSEs) and align the corresponding regions between the two structures. One way to do this is by using the DSSP (Define Secondary Structure of Proteins) online tool (<http://bioinformatica.isa.cnr.it/SUSAN/NAR2/dsspweb.html>), an application that assigns secondary structure to proteins, and a script in Visual Studio Code utilizing the biopython library in tandem. After generating the dssp files from their corresponding pdb files, we run the following code:

```
1 from Bio.PDB import PDBParser, Superimposer
2
3 # PDB file paths
4 pdb1_path = "7e2y.pdb"
5 pdb2_path = "AF-D2K8P9-F1-model_v4.pdb"
6
7 # DSSP file paths
8 dssp1_path = "7e2y.dssp"
9 dssp2_path = "AF-D2K8P9-F1-model_v4.dssp"
10
11 # Parse the PDB files to get the structures
12 pdb_parser = PDBParser()
13 structure1 = pdb_parser.get_structure("7e2y", pdb1_path)
14 structure2 = pdb_parser.get_structure("AF-D2K8P9-F1-model_v4", pdb2_path)
15
16 # Parse the DSSP files to get the secondary structure elements
17 def parse_dssp(dssp_path):
18     with open(dssp_path) as f:
19         lines = f.readlines()[28:] # DSSP secondary structure starts from
20         line 29
21         return [line[16] for line in lines if len(line) > 16 and line[16] in '
22         EHBIGST']
23
24 dssp1 = parse_dssp(dssp1_path)
25 dssp2 = parse_dssp(dssp2_path)
26
27 # Get the atom coordinates for secondary structure elements
28 def get_secondary_structure_atoms(structure, dssp):
29     atoms = []
30     for model in structure:
31         for chain in model:
32             for i, residue in enumerate(chain):
33                 # Get the residue number
34                 res_num = residue.get_id()[1]
35
36                 # If we have a secondary structure element for this residue
37                 if res_num < len(dssp) and dssp[res_num] != '-':
38                     try:
39                         atoms.append(residue["CA"])
40                     except:
```

```
39         continue
40     return atoms
41
42 atoms1 = get_secondary_structure_atoms(structure1, dssp1)
43 atoms2 = get_secondary_structure_atoms(structure2, dssp2)
44
45 # Modify this section to use the smaller number of atoms from atoms1 and
    atoms2
46 min_length = min(len(atoms1), len(atoms2))
47 super_imposer = Superimposer()
48 super_imposer.set_atoms(atoms1[:min_length], atoms2[:min_length])
49
50 print("secondary structure elements RMSD:", super_imposer.rms, "Angstrom")
```

In the beginning, the script imports necessary modules and defines file paths for the pdb and the dssp files. The pdb files contain the 3D atomic coordinates of the proteins, while the dssp files contain information about the secondary structure elements (alpha helices, beta sheets, etc.) of the proteins. The PDBParser object is used to parse the pdb files and get the protein structures. Next, the dssp files are parsed using the *'parse\_dssp'* function to get the secondary structure elements. This function reads each line of the dssp file starting from line 29 (as the actual dssp data begins from this line) and stores the secondary structure code (at position 16 in each line) if it corresponds to a known secondary structure element.

Subsequently, the *'get\_secondary\_structure\_atoms'* function is used to extract only the alpha carbon ('CA') atoms corresponding to secondary structure elements from each protein structure. It does so by iterating through each model, chain, and residue of the structure, and checking if the residue number corresponds to a secondary structure element in the dssp data.

Finally, the script calculates and prints the RMSD of the secondary structure elements between the two structures using the *'Superimposer'* class from BioPython. The RMSD is computed only for the smaller number of atoms between the two structures. The RMSD value gives a quantitative measure of the similarity between the secondary structure elements of the two protein structures.

The output in the VSCode terminal is: **secondary structure elements RMSD: 34.215 Angstrom**, which suggests significant conformational differences in the secondary structures, but this doesn't negate the potential for similar functionality. The comparison is subject to possible errors and uncertainties, since it was conducted between predicted and experimentally determined structures whose secondary structure elements do not align properly.