

## 1 Motivation

Proteins are cellular workhorses – they perform most of biological functions. Often, protein function includes uptake and release of other small molecules such as oxygen. These enter into the protein via "channels" or "pathways" – connected voids in the tightly packed protein structure. Examples of such proteins include hemoglobin that carries oxygen around our bodies, or myoglobin that stores oxygen in muscle cells (it is oxygen-rich myoglobin that makes fresh meat look red). Being able to determine voids and pathways inside proteins is important from both pragmatic (medical) and fundamental (how things work?) reasons. For example, a single mutation (error) in the oxygen carrying protein can change its ability to store and release oxygen resulting in *sickle cell anemia* disease. Computational methods that help understand biological function of proteins based on their structures speed up the processes of finding treatment for human disease. In fact, "bio computation" is used extensively in modern "drug discovery" – design of new medicines.

## 2 The assignment

Design, implement, and test a basic, simple algorithm that searches out for empty cavities in proteins. It does not have to be the most efficient or elegant one. It just has to work.

## 3 Definitions

### 3.1 Protein

A protein can be uniquely specified by  $(X,Y,Z)$  coordinates of its every atom  $i$ , see any file from the Protein Data Bank (PDB). Each atom is assumed to be a solid sphere of radius  $\rho(x_i, y_i, z_i)$ . Typically,  $1 \leq \rho \leq 2$  (in atomic length units called "Angstroms"). Proteins are tightly packed globs of atomic spheres, Fig. 1. Neighboring atom spheres can touch and even overlap. In fact, it is this intuitive and physically sound representation – CPK – that is utilized by most freely available visualization codes (e.g. `rasmol` or `pymol`) used to visualize protein structures. For example, carbon atoms are often displayed as a gray spheres, nitrogen are shown as smaller blue spheres. An atom whose name begins with an "M" would be displayed as a purple or magenta sphere. A typical protein contains between 500 to 5000 atoms.

### 3.2 Voids and Cavities

A point in space is said to belong to *void* space if a *probe sphere* of specified radius  $\rho_w$  can be placed at that point such that the sphere does not overlap with any protein atom. In biology, the relevant  $\rho_w = 1.4$  (radius of water molecule). A *cavity* is void space completely enclosed in the protein. That is there is no path of void points that connect any cavity point to the space *outside* the protein.

## 4 Specific tasks

Design, test, and implement an algorithm that finds all cavities in the input protein structure. Each cavity point needs to be identified within the accuracy of  $0.25$  ( $\text{\AA}$ ). The code will report a discrete set of spheres representing all of the cavities found in the protein. The spheres that belong to the same connected cavity should be placed no more than  $0.25$  units from each other along x, y, and z. In other words, you are representing each cavity with a set of spherical "pixels" at  $0.25$  resolution. It is OK to miss a cavity, if one needs higher than  $0.25$  resolution to find it. For example, suppose the probe radius  $\rho_w = 1.4$ , and the largest sphere that can fit anywhere inside the protein without touching any of the protein atoms (and not connected to the outside) has radius  $1.5$ . This cavity may be missed at  $0.25$  resolution, but will be identified at  $0.05$  resolution. However, if the largest sphere that can fit has radius of  $1.7$ , it must be found at  $0.25$  resolution, and a sphere of radius  $1.4$  must be placed in its center  $\pm 0.25$ .

You can also place up to additional 6 spheres around it, distance  $\pm 0.25$  along x, y, and z, if none of them touch any of the atoms. The cavity will be represented by this set of spheres. Due to finite resolution the exact number of spheres may differ (which is OK), as long as the cavity is found and represented by the appropriate spheres. Likewise, it is OK to report a false cavity if its elimination requires higher than 0.25 resolution, for example if higher resolution is required to determine that the putative cavity is actually connected to the outside.

#### 4.1 Formats

Your code assumes a very specific input format, called "PDB" (Protein Data Bank). We will use a simplified form of it shown in the example input and output files that you should download from the class site. The files are just the format examples, and do not necessarily have all the cavities present. The example is a cooked-up "spherical" protein. The output format has to be exactly the same as in the example. It is extremely important that all the spacings are identical to those in the examples. Field 2 is atom's number  $i$ . The first letter of field 3 is the name of the chemical element of that atom, e.g. "S" for sulfur (visualized as a yellow sphere). Field 4 is a 3-letter code for amino-acid name, field 5 is its number. You do not have to worry about those. Fields 6,7,8 are  $x_i, y_i, z_i$  coordinates of atom  $i$  of the protein or the sphere representing a cavity. Feel free to put 0.00 in the last but one field. This number is irrelevant. The last field is  $\rho_i$ . It varies from atom to atom: a typical value for carbon is 1.7, hydrogen 1.2, etc. All spheres representing cavity points must have the same  $\rho = \rho_w$ . The index for cavity spheres, field 2, must start at 10001 to distinguish them from the protein atoms. The index in field 5 must start at 5001. "CAV" in the 4th field indicates that the given sphere belongs to a cavity. Also, in the output file, the cavity records are always appended to the end of the input protein file, see example output.

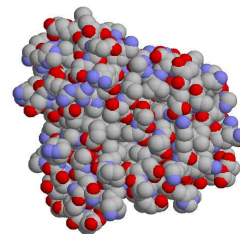


Figure 1: A typical small protein (myoglobin) in the so-called "CPK" or "space-fill" representation. Each atom is shown as a sphere of appropriate radius. Visualized by freely available code `rasmol`. Carbon atoms are gray spheres, oxygens are red, nitrogens are blue.

#### 4.2 The code

Strictly follow the programming guidelines, see below. C, C++, or PERL is preferred, but JAVA is OK too. Along with a project report, you will submit your code (If you must use Java: submit `.jar` in addition to the source). Ample comments are a must. Do not optimize the code for speed, but make sure it produces an answer within half an hour at most, for structures with up to 5000 atoms. It should give correct answer for a few, relatively simple test structures that we will use to test it (but you must come up with your own test cases, see below). We will test your code first for  $\rho_w = 1.4$ , but it should also give meaningful answers for other values of probe  $\rho$  within  $1 \leq \rho \leq 5$  interval. As well as for lower resolution such as 0.5. The code takes only three input parameters, and outputs just one file in the exact format as the example output you have. `mycode -i inputfile -o outputfile -probe 1.7 -resolution 0.25`

#### 4.3 What to report. Project stages

1. Stage 1. 20 points. During the first stage, each group gives a 3 min. in-class presentation to outline their plan. Followed by a 3 min. discussion with the rest of the class. This is your chance to get a feedback to see if you are on a right track. Time limit strictly enforced. Focus on general strategy and heuristics. Show results of your brainstorming the problem. No algorithm or implementation details are expected at this stage.
2. Stage 2. 50 points total. A few weeks later, each groups submits a typed progress report, up to 2.5 pages. Plus a zip file with the test structures, see below. The report MUST describe (with pictures) the algorithm you are planning to implement, ideas on how to test it, and heuristics you will have used (30 points). You must describe exactly, in detail, how you will distinguish cavity from void. You must also prepare, and submit (10 points), at least 3 test semi-realistic structures that you will use later (in stage 3) to test your code. Remember that proteins are compact "lumps" of atoms. Each structure should contain 1000 to 5,000

atoms, strictly following the correct input format. One structure, `solid.pdb`, should have no voids, the second, `cavity.pdb`, that you can easily make from `solid.pdb`, should have a single, fairly large cavity (to fit at least 10 distinct probe spheres at the required resolution), and the third one, `void.pdb`, should be produced from the second one, by boring a "tunnel" that connects the cavity to the outside, thus making the cavity a void. Note: for the actual testing in stage 3 you may need more of those test cases, but at this stage submit only the above three. If you already have some code at that stage (not mandatory), give the results on the above test, which will confirm the algorithm. But do not submit the code. Finally (10 points), attempt some kind of a computational feasibility analysis, along the lines we did in class. Focus on the most computationally intense part of your algorithm.

3. Stage 3. 80 points. Final typed report. By the end of class. No exceptions. Exact due date to be announced. No late reports under any circumstances. The report, up to 3 pages long must include a careful description of the algorithm (with schematics and diagrams), an outline of the implementation, and test results. Pictures are a must. You also submit your code, which we will compile and test. If the code does not pass our own simple test cases, including if we can not visualize your output because you did not adhere to the format specs, you are going to lose a lot of points.

**Submission** Strictly follow the "General Assignment Guidelines" (Group assignment) on the course web-site. Each group submits one PDF document, with roles of each partner clearly indicated. Everyone in the group will get the same score, assuming roughly equal amount of effort by group members. Otherwise, if the apparent amount of effort is highly unequal, the individual scores will be scaled accordingly.

**Programming guidelines** Follow "Programming guidelines" (see class website). In particular:

1. Your program should be submitted as a single zip archive.
2. The archive must include a README with a clear compilation instructions. Those should be extremely simple. You should assume that we will be testing your code on a very standard platform, e.g. UNIX CentOS 6. Do not use anything fancy. The code must compile!
3. The name of the archive should be: `last-name-of-student-1.last-name-of-student-N.ProjectCS2104.zip`
4. In the header of each class, you are required to include the following information about the assignment: (a) Name of IDE or compiler used. Use very standard ones. e.g. `g++` (b) Full names of all partners.
5. Use only standard libraries in your solutions.

## 5 Extra Credit (40 points)

All modern computers are multi-core. Parallelize the code you have developed above for your desktop/laptop to achieve a  $> 1.8$  speed-up relative to the single core. Note, that this is a non-trivial task, which will require changes to the algorithm. You can not simply assign various spatial regions to different CPUs since the voids may occur at the region boundary. A successful solution will always achieve  $> 1.8$  speed-up on  $N \geq 2$  CPUs, while producing *exactly the same (within machine precision)* answer as the single core version for any test case, including the trivial case when the input structure contains no voids (tightly packed atoms).

### 5.1 What to submit.

Two pages total including: (1) A description of the modifications to the original algorithm, with pictures. (2) Test cases, showing single core vs. multi-core results side-by-side. Include the "no voids" results. (3) Speed comparison of single core vs. multi-core versions for all the above test cases. Do not submit the source code, but be ready to send it upon request. Specify details of your computer and the parallel implementation (e.g. OpenMPI).