

Extraction of Robust Voids and Pockets in Proteins

Raghavendra Sridharamurthy¹, Harish Doraiswamy², Siddharth Patel¹, Raghavan Varadarajan¹ and Vijay Natarajan¹

¹Indian Institute of Science, Bangalore
²Polytechnic Institute of New York University, USA

Abstract

Voids and pockets in a protein refer to empty spaces that are enclosed by the protein molecule. Existing methods to compute, measure, and visualize the voids and pockets in a protein molecule are sensitive to inaccuracies in the empirically determined atomic radii. This paper presents a topological framework that enables robust computation and visualization of these structures. Given a fixed set of atoms, voids and pockets are represented as subsets of the weighted Delaunay triangulation of atom centers. A novel notion of (ϵ, π) -stable voids helps identify voids that are stable even after perturbing the atom radii by a small value. An efficient method is described to compute these stable voids for a given input pair of values (ϵ, π) .

Categories and Subject Descriptors (according to ACM CCS): I.3.6 [Computer Graphics]: Computer Graphics—Methodology and Techniques

1. Introduction

Protein molecules have a well packed structure, yet they contain cavities. A cavity refers to both voids (without openings) and pockets (with openings). These cavities play a key role in determining the stability and function of proteins.

Several methods have been proposed to locate such cavities in protein molecules. We focus our attention on geometric methods. Edelsbrunner et al. [EF94, EK05] and Liang et al. [LEF*98, LEW98] propose a definition that is based on the theory of alpha shapes and discrete flows in Delaunay triangulations. Kim et al. [KCS*10, KS12] propose a definition of cavities based on an alternate representation of a set of atoms called beta shapes that faithfully captures proximity. Tools based on the above approach are available and widely used [DOT*06, KLE04, KRSC12]. Till and Ullmann [TU10] employ a graph theoretic algorithm to identify cavities and compute their volume. Parulek et al. [PTRV12] use graph based methods on the implicit representation of molecular surfaces to identify pockets and potential binding sites. Varadarajan et al. [CBV02] employ a Monte Carlo procedure to position water molecules together with a Voronoi region-based method to locate empty space. They discuss the importance of accurate identification of cavities for the study of protein structure and stability. Novel Voronoi diagram-based techniques for the extraction and visualization of cavities have also been developed from the viewpoint of studying and interactively exploring access

paths to active sites [POB*06, PKKO07, LBH11, LBBH12]. Krone et al. [KFR*11] present a visualization tool for interactive exploration of protein cavities in dynamic data.

1.1. Motivation and Problem Statement

The input used for calculations in previous work come from x-ray crystallography data or other lower resolution data. Previous cavity detection methods are sensitive to inaccuracies that are inherent in the crystallographic measurements. While the measurements may guarantee high resolution, it is important to note that even small inaccuracies may cause a significant difference in the reported number of cavities. The inaccuracies may also arise due to some fundamental limitations such as the notion of radii of atoms, which is determined empirically. For example, as illustrated in Figure 1, presence of such inaccuracies may result in a cavity detection method to report two distinct but large cavities in place of one or report very small volume cavities. Figure 2 illustrates the problem as it occurs in a lysozyme protein.

We aim to develop an interactive method to compute robust cavities in proteins. We achieve this by enabling the user to reduce, if not completely eliminate, the inaccuracies mentioned earlier. We define what robustness would mean in this context, why such a notion is important and also demonstrate the robustness of the method. As an auxiliary task, we visu-



Figure 1: *Left:* Two cavities that are apparently very near to each other may be a single cavity. *Right:* A very small cavity may be reported whereas no such cavity may exist.

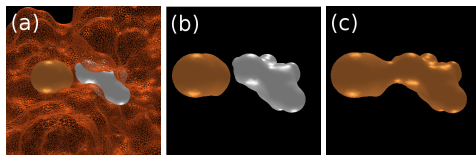


Figure 2: (a) Two voids that appear very near to each other in a lysozyme protein (PDB ID: 200l). The solid surface represents the voids while the wireframe represents the molecule. (b) Zoomed-in view of the voids without the molecular skin surface. (c) The two voids may be a single void.

alize the protein molecules and their cavities in an interactive manner and compute their properties.

From the biologist’s point of view, obtaining a stable protein is the starting point of many applications, from in-vitro studies of binding and interactions, to using the protein as an antigen or vaccine. Whereas surface pockets often form part of the active site of enzymes or interacting sites for other proteins, internal voids are often relevant structurally as features that affect the overall stability of the protein. It is established that filling up internal voids improves the packing of the protein thus increasing stability. In this respect, detecting and visualizing structurally robust cavities inside the protein informs the biologist on which mutations to perform to improve internal packing and get a stable protein.

1.2. Results

The main results described in this paper are in the context of a novel definition and method for computing robust and stable voids in proteins. We employ a simple and succinct structure called the alpha complex to represent protein molecules. The alpha complex is a simplicial complex that can be stored as a filtration, a series of simplicial complexes K^i with $K^i \subset K^{i+1}$. With the aim of computing a set of voids that are stable with respect to small perturbations in the atom radii, we develop a method that modifies the radii of a select set of atoms symbolically by systematically processing and modifying the filtration. We show that this modification results in controlled changes in the number and properties of voids and does not violate key properties of the filtration. The method also supports the elimination of very small or insignificant voids as measured by the notion of topological persistence [ELZ02]. We also develop software to visualize the stable cavities together with the molecule, and to calculate cavity volume and surface areas. An extended technical report presents additional experimental results [SDF*13].

2. Background

In this section, we briefly introduce the topological background required to define and represent the structure of biomolecules [Mun84, Ede10, Ede04].

Simplicial Complex. A k -simplex σ is the convex hull of $k + 1$ affinely independent points. A vertex, edge, triangle, and tetrahedron are k -simplices of dimension $0 - 3$. A simplex τ is a *face* of σ , $\tau \leq \sigma$, if it is the convex hull of a non-empty subset of the $k + 1$ points. A *simplicial complex* K is used to represent a topological space and is a finite collection of simplices such that (a) $\sigma \in K$ and $\tau \leq \sigma$ implies $\tau \in K$, and (b) $\sigma_1, \sigma_2 \in K$ implies $\sigma_1 \cap \sigma_2$ is either empty or a face of both σ_1 and σ_2 . A *subcomplex* of K is a simplicial complex $L \subseteq K$.

Voronoi diagram and Delaunay triangulation. Let $S \subseteq \mathbb{R}^d$ be a finite set of points. The *Voronoi cell* V_p , of a point $p \in S$, is the set of points in \mathbb{R}^d whose Euclidean distance to p is smaller than or equal to any other point in S . The collection of Voronoi cells of all points in S partitions \mathbb{R}^d , and is called the *Voronoi diagram* (Figure 3(a)). The *Delaunay triangulation* D of S is the dual of the Voronoi diagram and partitions the convex hull of S , see Figure 3(b). The above definitions can be extended to a set of balls or weighted points by choosing an appropriate measure of distance between a weighted point p and a point in \mathbb{R}^d . The *power distance* between p and a point $x \in \mathbb{R}^d$ is equal to $\pi_p(x) = \|x - p\|^2 - w_p$, where w_p is the weight of p .

Alpha Complex. Molecules are often represented using a space-filling model such as a union of balls. The weighted Voronoi diagram may be extended to represent the contribution from each atom to the union of balls. Consider an atom p . Define B_p as an open ball having the radius of the atom p . Let V_p be the weighted Voronoi cell corresponding to p , where the weight is equal to the square of the atom radius. The contribution from each atom p is equal to $B_p \cap V_p$, the intersection between the ball corresponding to the atom and the weighted Voronoi cell of p . The corresponding dual structure is a subcomplex of the weighted Delaunay triangulation and called the *dual complex*, see Figure 4.

Edelsbrunner et al. [EKS83, EM94, Ede92] consider a growth model where the ball radii grow, and track the changes in the dual complex. The growth parameter, α , corresponds to a radius $\sqrt{r_p^2 + \alpha^2}$ for a ball centered at p with radius r_p . The weight of the point $w(p)$ increases to $w(p) + \alpha^2$. Note that $\alpha = 0$ corresponds to no growth. The dual complex corresponding to a set of balls after they are grown by α is called the *alpha complex*.

Given a simplicial complex K , a finite sequence, $\emptyset = K^0, K^1, \dots, K^m = K$, of subcomplexes of K is a *filtration* if $K^0 \subset K^1 \subset \dots \subset K^m$. The *rank* of a subcomplex refers to its position in the filtration. The set of alpha complexes obtained by varying α from $-\infty$ to ∞ is a filtration of the Delaunay triangulation. In particular, we consider the filtration that is

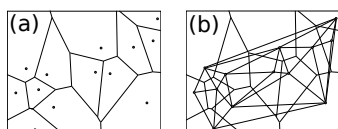


Figure 3: (a) Voronoi diagram of a point set in \mathbb{R}^2 . (b) The Delaunay complex is the dual of the Voronoi diagram.

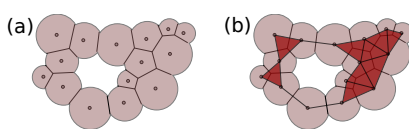


Figure 4: (a) Intersection of the weighted Voronoi diagram and the union of balls. (b) The dual complex is the dual of this partition of the union of balls that captures the incidence relationship.

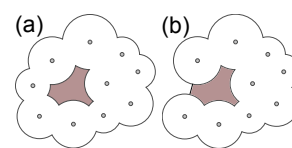


Figure 5: (a) Void and (b) Pocket in a collection of 2D balls.

generated by inserting the simplices one at a time and ties broken based on the dimension of the simplex.

Voids and Pockets. Let the alpha complex K represent a molecule at a given value α and D be the Delaunay complex of the weighted point set. A *cavity* is a maximally connected component of the complement $D - K$. *Voids* and *pockets* are cavities that are, respectively, bounded and not bounded by the union of balls. Figure 5 illustrates a void and a pocket in 2D. The alpha complex helps represent and track voids via the growth process. *Topological persistence* measures the lifetime of a void [ELZ02] and is equal to the difference between the ranks of alpha complexes in the filtration when the void is created and destroyed.

3. Robust Voids and their Computation

We introduce a notion of robust voids based on two parameters, one local and another global. The local parameter is referred to as stability and the global parameter is specified by topological persistence. In order to simplify the description, we assume that the voids are computed for the α -complex corresponding to $\alpha = 0$. However, the definitions, methods, and subsequent analysis are valid for all values of α .

3.1. ϵ -stable and π -persistent voids

Consider the interval $[-\epsilon, \epsilon]$ of α values, where $\epsilon \geq 0$. A void is called an ϵ -stable void if it remains a single connected void within all α -complexes for α values in the range $[-\epsilon, \epsilon]$. In other words, using the lifetime terminology, the void is born, possibly split into multiple components, and destroyed at α -values that lie strictly outside of this interval. A void is π -persistent if its topological persistence is greater than π i.e., the void size measured in terms of its lifetime is greater than π . Combining the two notions of robustness, we call a void to be (ϵ, π) -stable if it is both ϵ -stable and π -persistent.

The above definitions help measure the stability of the voids when the radii are perturbed by a small value. The local parameter considers perturbation within a small interval centered at the α -value of interest whereas the global parameter measures the size of the void in terms of its lifetime in the filtration. Voids of interest may often not be stable with respect to both notions. For example, a large sized void (π -persistent for some large π) may be born within the interval $[-\epsilon, \epsilon]$. However, note that a small perturbation in the radii

of atoms that line the surface of the void could result in an earlier birth time, hence making the void to be ϵ -stable. We aim to extract all voids that are either stable as is or can be made stable via a small perturbation.

3.2. Computing (ϵ, π) -voids

The location of the atoms that constitute a protein molecule together with their van der Waals radii is obtained from the protein data bank in pdb format. Given ϵ and π , we compute the set of (ϵ, π) -stable voids as follows.

1. Compute the weighted Delaunay triangulation of the input [KLE04]. The atom centers form the set of points that are weighted using their van der Waals radii.
2. Build the alpha shape spectrum [Ede92], which is a filtration of the weighted Delaunay triangulation.
3. **Modify the filtration based on the value of ϵ .**
4. Compute the set of (ϵ, π) -stable voids by identifying all voids [LEF*98] of the modified filtration at $\alpha = 0$, and retaining only those voids that have persistence at least π .

Modifying the filtration. The filtration of the weighted Delaunay triangulation as defined by the α -values provides an explicit representation of the birth/death times of each void and the evolution during its lifetime. We propose to alter the birth/death times of voids by modifying the filtration instead of directly modifying radii of atoms that line the surface of the void. While the latter approach follows directly from the definition, it is cumbersome and computationally inefficient. For example, varying the radii without explicit control may lead to changes in the triangulation and the alpha complex. These changes need to be explicitly tracked, else they may lead to inconsistencies between the alpha complex that represents the molecule and the space-fill model. Resolving such inconsistencies would necessitate the re-computation of all representations. On the other hand, the former approach is simpler and computationally efficient.

Delayed simplex insertion. One or more simplices are inserted to obtain a rank $i + 1$ simplicial complex from a rank i simplicial complex in the filtration. Higher ranks correspond to higher values of α . The topology of voids may change when the simplices are inserted. In particular, the insertion of a triangle may either not affect any void, split a void into two, or create a new void. On the other hand, the insertion of a tetrahedron always destroys a void. These topology changes

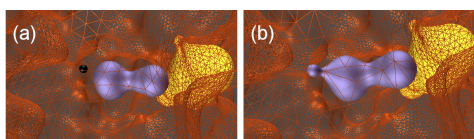


Figure 6: Robust voids in protein 2CI2. (a) Voids at $\alpha = 0$ (b) $(0.3, 0.01)$ -stable voids. The violet and black voids merge to form a single stable void.

may therefore be avoided by delaying the insertion of the simplices that change the topology of voids. Let K^j be the alpha complex corresponding to $\alpha = \epsilon$. Consider the set of simplices, Σ , inserted into the filtration for values of α in the range $[-\epsilon, \epsilon]$. Let $\Sigma_t \subset \Sigma$ be the set of the triangles that split a void, and $\Sigma_T \subset \Sigma$ be the set of tetrahedra. We delay the insertion of simplices σ_i in Σ_t and Σ_T such that $\sigma_i \notin K^j$ but $\sigma_i \in K^l$, where $K^j \subset K^l \subset D$. Simplicial complexes in the filtration of the weighted Delaunay triangulation and the order of simplices that are inserted to generate the filtration satisfy several containment and incidence properties. These properties should be satisfied for the modified filtration as well. Towards this, we propose a conservative but computationally efficient approach to modify the filtration:

1. Move all tetrahedra in Σ_T to the end of the filtration. All such tetrahedra are present in D but not in any $K^i \subset D$.
2. For each triangle in Σ_t , find its incident tetrahedra τ_1, τ_2 .
3. Delay the insertion of the triangle and the two tetrahedra, τ_1 and τ_2 , to the end of the filtration.

Implication of the modified filtration. Consider the delayed insertion of a triangle that causes a void to split. The void is no longer split into two and instead reported as a single connected ϵ -stable void. The delay corresponds to shrinking the atoms centered at the vertices of the triangle. However, note that the radii are not yet modified. We optionally modify the radii later for further analysis of the void. A triangle that creates a void is left untouched and the corresponding void is also declared to be ϵ -stable. The triangle insertion may be advanced to ensure that the void is created outside the interval. This corresponds to a small increase in the radii of the atoms centered at the vertices of the triangle. We choose not to explicitly advance the triangle insertion because it does not affect the results for small values of ϵ . After the filtration is modified as described above, we recompute the voids from the alpha complex.

We compute the persistence of the ϵ -stable voids obtained and retain only those having persistence greater than π . This pruned set of voids are (ϵ, π) -stable. Note that the persistence is computed with respect to the original filtration. The above notion of stability can be extended to pockets as well.

Analysis. Let m be the number of simplices in the Delaunay triangulation of the input protein having n atoms, $m = O(n^2)$. Computing the set of voids takes $O(m\alpha(m))$ time using the union-find data structure. Here, α is the inverse Ackermann

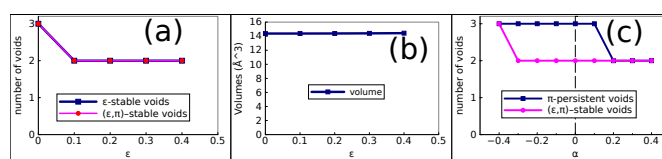


Figure 7: Properties of stable voids of the protein 2CI2. (a) Number of voids decreases for larger ϵ . (b) Total volume of voids increases only marginally. (c) (ϵ, π) -stable voids in the neighbourhood of $\alpha = 0$.

function. Given ϵ , Σ_t and Σ_T is computed in $O(m)$ time using a sequential search over the filtration. Identifying the set of tetrahedra incident on triangles in Σ_t , and moving all the simplices to the end of the filtration takes $O(m)$ time. Thus the time required to modify the filtration is $O(m\alpha(m))$.

4. Results

Our software *RobustVoids* computes the alpha complex and supports its visualization for different α -values. The set of (ϵ, π) -stable cavities are computed using the values of ϵ and π specified by the user. The software also reports the cavity volume and surface area.

Visualization of stable cavities. Figure 6(a) shows protein 2CI2, which has 3 voids. Using a value of $\epsilon = 0.3$ and $\pi = 0.01$ results in two $(0.3, 0.01)$ -stable voids, see Figure 6(b). A value of $\epsilon = 0.3$ is equivalent to an increase / decrease of the atom radius by at most 0.33\AA , which is within the tolerance limit of 0.5\AA . Modifying the filtration and computing the stable voids for this protein takes 0.1 s.

Properties of stable cavities. Figure 7(a) and 7(b) plots the number and volume of (ϵ, π) -stable voids for increasing values of ϵ . Increasing the value of ϵ implies that voids from a larger α range are considered. This could potentially increase the number of ϵ -stable voids. However, such voids usually have low persistence and are not (ϵ, π) -stable. We use a constant value of $\pi = 0.01$ in all experiments. The total volume of all stable voids increases marginally ($< 1\%$) with increasing ϵ . The merging of two nearby voids into a single stable void does not effect the total volume. However, volumes of individual voids could change drastically. The volume of the stable void in Figure 6(b) is approximately equal to the sum of the volumes of the unmerged voids. The volumes are verified against known results [CBV02]. Figure 7(c) shows that the number of $(0.3, 0.01)$ -stable voids do not change for $\alpha \in [-0.3, 0.3]$ as compared to the number of 0.01-persistent voids.

Summary. We have defined a novel notion of robust cavities that is insensitive to the perturbation of the atomic radii. Robust cavities are computed via a controlled modification of the filtration that represents the molecule and its cavities. This computation benefits a biologist who can now perform tedious mutation-based experiments only on these cavities.

Acknowledgements. This work was supported by a grant from Department of Science and Technology, India (SR/S3/EECE/0086/2012).

References

- [CBV02] CHAKRAVARTY S., BHINGE A., VARADARAJAN R.: A procedure for detection and quantitation of cavity volumes in proteins. *Journal of Biological Chemistry* 277, 35 (2002), 31345–31353. 1, 4
- [DOT*06] DUNDAS J., OUYANG Z., TSENG J., BINKOWSKI A., TURPAZ Y., LIANG J.: CASTp: computed atlas of surface topography of proteins with structural and topographical mapping of functionally annotated residues. *Nucleic acids research* 34, 2 (2006), W116–W118. 1
- [Ede92] EDELSBRUNNER H.: *Weighted alpha shapes*. University of Illinois at Urbana-Champaign, Department of Computer Science, 1992. 2, 3
- [Ede04] EDELSBRUNNER H.: Biological applications of computational topology. In *Handbook of Discrete and Computational Geometry*, Goodman J. E., O'Rourke J., (Eds.). CRC Press, 2004, pp. 1395–1412. 2
- [Ede10] EDELSBRUNNER H.: *Computational Topology. An Introduction*. Amer. Math. Soc., 2010. 2
- [EF94] EDELSBRUNNER H., FU P.: *Measuring space filling diagrams and voids*. Tech. rep., UIUC-BI-MB-94-01, Beckman Inst., Univ. Illinois, Urbana, Illinois, 1994. 1
- [EK05] EDELSBRUNNER H., KOEHL P.: The geometry of biomolecular solvation. *Combinatorial & Computational Geometry* 52 (2005), 243–275. 1
- [EKS83] EDELSBRUNNER H., KIRKPATRICK D., SEIDEL R.: On the shape of a set of points in the plane. *IEEE Transactions on Information Theory* 29, 4 (1983), 551–559. 2
- [ELZ02] EDELSBRUNNER H., LETSCHER D., ZOMORODIAN A.: Topological persistence and simplification. *Discrete & Computational Geometry* 28, 4 (2002), 511–533. 2, 3
- [EM94] EDELSBRUNNER H., MÜCKE E.: Three-dimensional alpha shapes. *ACM Transactions on Graphics (TOG)* 13, 1 (1994), 43–72. 2
- [KCS*10] KIM D.-S., CHO Y., SUGIHARA K., RYU J., KIM D.: Three-dimensional beta-shapes and beta-complexes via quasi-triangulation. *Computer-Aided Design* 42, 10 (2010), 911–929. 1
- [KFR*11] KRONE M., FALK M., REHM S., PLEISS J., ERTL T.: Interactive exploration of protein cavities. In *Computer Graphics Forum* (2011), vol. 30, pp. 673–682. 1
- [KLE04] KOEHL P., LEVITT M., EDELSBRUNNER H.: Proshape: understanding the shape of protein structures. *Software at biogeometry.duke.edu/software/proshape* (2004). 1, 3
- [KRSC12] KIM D.-S., RYU J., SHIN H., CHO Y.: Beta-decomposition for the volume and area of the union of three-dimensional balls and their offsets. *Journal of Computational Chemistry* (2012). 1
- [KS12] KIM D.-S., SUGIHARA K.: Tunnels and voids in molecules via voronoi diagram. In *Proc. Symp. Voronoi Diagrams in Science and Engineering (ISVD)* (2012), pp. 138–143. 1
- [LBBH12] LINDOW N., BAUM D., BONDAR A., HEGE H.: Dynamic channels in biomolecular systems: Path analysis and visualization. In *Proc. IEEE Symposium on Biological Data Visualization (BioVis)* (2012), pp. 99–106. 1
- [LBH11] LINDOW N., BAUM D., HEGE H.: Voronoi-based extraction and visualization of molecular paths. *Visualization and Computer Graphics, IEEE Transactions on* 17, 12 (2011), 2025–2034. 1
- [LEF*98] LIANG J., EDELSBRUNNER H., FU P., SUDHAKAR P., SUBRAMANIAM S.: Analytical shape computation of macromolecules: II. inaccessible cavities in proteins. *Proteins Structure Function and Genetics* 33, 1 (1998), 18–29. 1, 3
- [LEW98] LIANG J., EDELSBRUNNER H., WOODWARD C.: Anatomy of protein pockets and cavities. *Protein Science* 7, 9 (1998), 1884–1897. 1
- [Mun84] MUNKRES J.: *Elements of Algebraic Topology*, vol. 2. Addison-Wesley Menlo Park, CA, 1984. 2
- [PKK07] PETŘEK M., KOŠINOVÁ P., KOČA J., OTYEPKA M.: MOLE: A Voronoi diagram-based explorer of molecular channels, pores, and tunnels. *Structure* 15, 11 (2007), 1357–1363. 1
- [POB*06] PETŘEK M., OTYEPKA M., BANÁŠ P., KOŠINOVÁ P., KOČA J., DAMBORSKÝ J.: Caver: a new tool to explore routes from protein clefts, pockets and cavities. *BMC bioinformatics* 7, 1 (2006), 316. 1
- [PTRV12] PARULEK J., TURKAY C., REUTER N., VIOLA I.: Implicit surfaces for interactive graph based cavity analysis of molecular simulations. In *Biological Data Visualization (BioVis), 2012 IEEE Symposium on* (2012), pp. 115–122. 1
- [SDP*13] SRIDHARAMURTHY R., DORAISWAMY H., PATEL S., VARADARAJAN R., NATARAJAN V.: *Extraction of robust voids and pockets in proteins*. Tech. Rep. IISC-CSA-TR-2013-3, Department of Computer Science and Automation, Indian Institute of Science, <http://www.csa.iisc.ernet.in/TR/2013/3/>, 2013. 2
- [TU10] TILL M. S., ULLMANN G. M.: Mcvol-a program for calculating protein volumes and identifying cavities by a monte carlo algorithm. *Journal of molecular modeling* 16, 3 (2010), 419–429. 1