

pClay: A Precise Parallel Algorithm for Comparing Molecular Surfaces

Georgi D. Georgiev*, Kevin F. Dodd* and Brian Y. Chen¹

Department of Computer Science and Engineering, Lehigh University, USA

Abstract

Comparing binding sites as geometric solids can reveal conserved features of protein structure that bind similar molecular fragments and varying features that select different partners. Due to the subtlety of these features, algorithmic efficiency and geometric precision are essential for comparison accuracy. For these reasons, this paper presents pClay, the first structure comparison algorithm to employ fine-grained parallelism to enhance both throughput and efficiency. We evaluated the parallel performance of pClay on both multicore workstation CPUs and a 61-core Xeon Phi, observing scaleable speedup in many thread configurations. Parallelism unlocked levels of precision that were not practical with existing methods. This precision has important applications, which we demonstrate: A statistical model of steric variations in binding cavities, trained with data at the level of precision typical of existing work, can overlook 46% of authentic steric influences on specificity ($p \leq .02$). The same model, trained with more precise data from pClay, overlooked 0% using the same standard of statistical significance. These results demonstrate how enhanced efficiency and precision can advance the detection of binding mechanisms that influence specificity.

2012 ACM Subject Classification Applied Computing → Molecular Structural Biology; Computing Methodologies → Volumetric Models; Parallel algorithms

Keywords and phrases Specificity Annotation, Structure Comparison, Cavity Analysis

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

1 Introduction

Molecular shape and electric fields have a strong influence on binding specificity. At binding interfaces, complementary molecular shapes can accommodate some ligands and hinder those that fit poorly. Electric fields attract molecules with complementing charges and repel others. This connection, between molecular recognition and the geometric complementarity of surfaces and fields, is evidence by which human investigators infer the roles of individual mechanisms in function. Comparison software can detect this kind of evidence and use it to make similar inferences. Some methods detect proteins with geometrically conserved binding sites, supporting the inference that they bind similar partners [11, 7, 4, 26, 32, 10, 14, 17]. Other methods find variations in the electric fields near binding sites, suggesting that they accommodate differently charged ligands [18, 5, 28, 33]. These techniques, and their potential for large scale and accurate applications, depend on rapid and precise digital representations of molecular shape, which are the focus of this paper.

Rapid and precise algorithms can integrate many observations to support inferences that are impossible with single comparisons. For example, a single comparison does not provide a frame of reference that would be needed to assess whether or not two binding sites are different enough that they have different binding preferences. After all, conformational variations and single mutations can occur in many ways that change nothing about binding. This kind of inference is traditionally reserved for experts with a wealth of biochemical experience. However, statistical models can be trained on the steric differences between closely related ligand binding sites that prefer the same ligands. In such cases, structures of close homologs or even single mutants could provide the primary data, but the subtle variations needed to train the model would have to be found with many individual comparisons. Once trained,

¹ Correspondence: chen@cse.lehigh.edu * Equal contribution



© Georgi D. Georgiev, Kevin F. Dodd and Brian Y. Chen;
licensed under Creative Commons License CC-BY

42nd Conference on Very Important Topics (CVIT 2016).

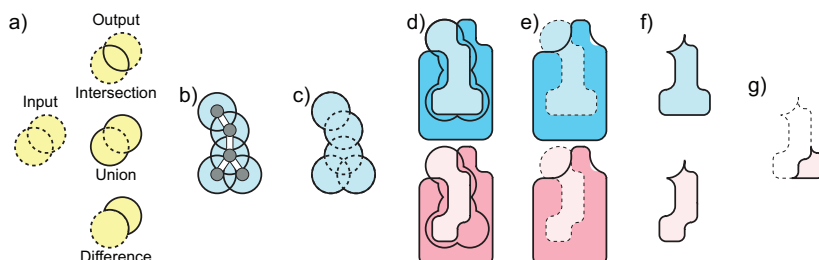
Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:12



Leibniz International Proceedings in Informatics

LIPIcs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

46 the statistical model provides a frame of reference that reveals steric variations that are too
 47 large to be typical of binding sites with the same binding preferences. The large variations
 48 found would therefore be indicators of binding sites that have different binding preferences
 49 [6]. To support and advance statistical models like these, this paper presents pClay, the first
 50 structure comparison algorithm that maximizes precision and computational throughput
 51 using arbitrary precision representations and parallel algorithms.



■ **Figure 1** CSG operations on Molecular Surfaces. a) Basic CSG operations. Input solids are yellow with dotted outlines. Outputs have solid outlines. b) Ligand with grey atoms and white bonds, with spheres centered on each atom (light blue). c) The union of atom-centered spheres. d) Two molecular surfaces (blue, red) in complex with two ligands shown as sphere unions (black lines). e) CSG difference of the sphere unions minus molecular surfaces (dotted lines), shown with molecular surfaces (blue and red, no outline) and envelope surfaces (black outline). f) Intersection of differences with envelope surfaces (light blue and red). g) The CSG difference between binding cavities reveals a variation in steric hindrance that causes differences in binding preferences.

52 pClay performs geometric comparisons using Constructive Solid Geometry (CSG) opera-
 53 tions (fig. 1a) on analytically represented three dimensional solids. These operations, which
 54 include unions, intersections and differences, can be combined like arithmetic operators to
 55 sculpt a geometric solid. This sculptural nature of CSG inspires both the name pClay, a
 56 portmanteau of “protein” and “clay,” and also the solid geometric approach to the analysis of
 57 protein shape that pClay makes possible. For example, the union of large spheres centered at
 58 ligand atoms can represent the neighborhood of a ligand (fig. 1b,c). The difference between
 59 the spheres and the molecular surface of a receptor can describe the solvent-accessible binding
 60 cavity in the receptor (fig. 1d,e). The CSG difference between one binding cavity and another
 61 is the cavity region that is solvent accessible in one protein and inaccessible in the other (fig.
 62 1g). This difference, the variation between the two cavities, could be small, when binding
 63 preferences are similar, or large, when steric hindrance creates differences in specificity.

64 The utility of these computations can be seen in multiple applications: When applying
 65 this approach to the S1 subsites of trypsins and elastases, we observed that it could identify
 66 threonine 226 which, in elastases, sterically hinders the longer substrates preferred by trypsins
 67 that might otherwise bind [8]. To illustrate the importance of precision, that region of
 68 hindrance is only 50 percent larger than a carbon atom (31 \AA^3). A similar approach identified
 69 “gatekeeper” residue 338 in the tyrosine kinases [13], which creates steric clash with larger
 70 drugs [21]. We have also observed that a CSG-based comparison of electrostatic isopotentials
 71 can reveal single amino acids crucial for selecting ligands in the in the cysteine proteases [5]
 72 and for stabilizing the three interfaces of the SMAD trimer [28]. Experimental validation has
 73 demonstrated the correctness of our prediction that arginine 235 forms critical electrostatic
 74 interactions for the activity of the ricin toxin [33]. By making CSG analysis possible on
 75 geometric solids that are exact, up to machine precision, pClay ensures that subtle but
 76 influential details cannot be overlooked.

77 The precision that pClay achieves derives from solids that have analytical representations,
 78 like spheres and tetrahedra. pClay can assemble these primitives into solvent excluded

79 regions, which we call *molecular solids*. The boundary of a molecular solid is the classic
80 molecular surface, also known as the solvent excluded surface, which was originally developed
81 by Richardson and others [20, 9]. While we can construct molecular solids with CSG
82 operations on many individual primitives, pClay exploits molecular properties to sidestep
83 those operations and achieve greater efficiency. The resulting molecular solids avoid the
84 “photocopier effect”, where multiple CSG operations can accumulate geometric errors. They
85 can also be exported as triangle meshes, generated at an arbitrary degree of precision, for
86 compatibility with other software.

87 pClay enhances computational efficiency through parallelism. We achieve parallelism in
88 pClay in a number of ways, most notably by recasting Marching Cubes, a traditional method
89 for implementing CSG operations [22, 16], into a series of parallel breadth first searches
90 (BFS). In pClay, we use BFS to traverse cubic lattices and identify contiguous regions of cubes
91 within defined boundary regions. These breadth first traversals can be distributed evenly
92 across arbitrary numbers of threads. By dividing the computation in this way, parallelism
93 can make comparisons faster and also enable more detail to be considered. This advancement
94 stands in qualitative contrast with existing efforts to parallelize structure comparisons (e.g.
95 [7]), where throughput was increased without enhancing precision. To demonstrate the
96 parallel scalability of our method, pClay was tested on both modern multicore processors as
97 well as on a Xeon Phi, a manycore processor with 61 cores.

98 Relative to existing methods, pClay is the first algorithm to use arbitrarily precise
99 representations of molecular surfaces for protein structure comparison. It is also the first
100 structure comparison method to use fine grained parallelization, enhancing both precision and
101 computational throughput. Several methods do employ arbitrarily precise representations of
102 the molecular surface, using NURBs [2], alpha shapes [31] or spherical coordinates [25, 27],
103 but they are used for visualization and have not been integrated into comparison algorithms.
104 Other methods parallelize structure comparison to refine representations of binding sites
105 [7], to accelerate database searches [19], or create cloud-based search services [15], but use
106 parallelism to enhance throughput and not also precision. To our knowledge, pClay is the
107 first integration of arbitrary precision and parallelism into a structure comparison method.

108 **2** Methods

109 As input, pClay accepts a collection of geometric solids and an expression of CSG operations.
110 We convert the CSG expression into a binary tree, a CSG tree, where the nodes of the
111 tree are geometric solids. The input solids, which include spheres, spindles, tetrahedra or
112 molecular surfaces, are leaves on the CSG tree, while the result of CSG operations are the
113 nonleaf nodes. The final result of all operations, the root node, is the output. pClay can also
114 generate a closed triangular mesh at user-defined resolutions to approximate the boundary
115 of the output.

116 To perform CSG operations, pClay implements a parallel version of Marching Cubes [22]
117 (Section 2.1), which we summarize below. Our method requires three basic functions to be
118 performed by every node in the CSG tree. These functions are *containsPoint()*, *intersectSeg-*
119 *ment()*, and *findSurfaceCubes()*. Given any point p in three dimensions, *containsPoint(p)*
120 determines exactly if p is inside or outside the solid. A point exactly on the surface is said
121 to be inside the solid. Second, given a line segment s , *intersectSegment(s)* determines all
122 points of intersection between the surface of the operand and s , as well as the interior or
123 exterior state of each interval on the segment. Finally, given a cubic lattice l that surrounds
124 the primitive, *findStartingCubes(l)* finds a few cubes of the lattice that are *surface cubes*,
125 having at least one corner inside and one corner outside the solid. These cubes are used

126 to initiate a parallel breadth first search for all surface cubes, called *findAllSurfaceCubes()*,
 127 which is implemented once for all primitives (Section 2.2). To implement leaf nodes it is
 128 thus sufficient to describe how these basic functions are implemented for that solid. Nonleaf
 129 nodes implement the basic functions as logical operations, as we will explain in Section 2.3.

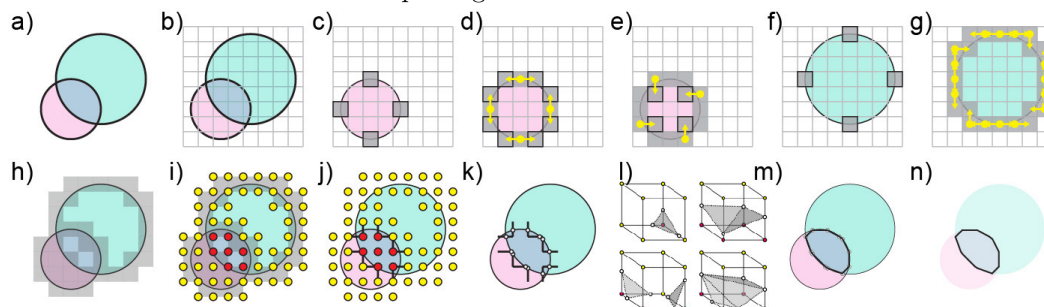
130 Below, we first describe how the output approximations are generated using a paral-
 131 lelization of Marching Cubes and how we find all surface cubes beginning the output from
 132 the starting cubes generated by the basic function. We next explain how the three basic
 133 functions are implemented for every primitive. Finally, we detail how the basic functions are
 134 implemented in nonleaf nodes.

135 2.1 Parallel Marching Cubes

136 As input, Marching Cubes accepts a set of geometric solids (fig. 2a), which we will refer to
 137 as operands, and a CSG expression tree to be performed on the operands. It also accepts a
 138 resolution parameter in angstrom units that specifies the degree to which the result of the
 139 CSG expression should be approximated in the output.

140 We begin by defining an axis aligned cubic lattice surrounding the input operands, where
 141 each cube has sides equal to the user-specified resolution parameter (fig. 2b). This step is
 142 performed by examining the sizes of all operands and the related CSG operations.

143 Once the lattice is defined, we invoke *findStartingCubes(l)* on each input solid (fig. 2c,f).
 144 The surface cubes identified are provided as input to *findAllSurfaceCubes()*, which identifies
 145 all remaining surface cubes of all inputs solids in parallel (fig. 2h). The process of identifying
 146 surface cubes for all input solids also necessarily determines the interior/exterior state of the
 147 points on these cubes in relation to specific solids. We then compute the interior/exterior
 148 state of these points in relation to all other solids in an embarrassingly parallel manner. Once
 149 this assessment is made for any point, we can access whether that point is inside or outside
 150 the output region (fig. 2i). In this way, we find the subset of cubes that contain a corner
 151 inside and a corner outside the output region.



■ **Figure 2** a) Input operands (red, green). b) Cubic lattice around operands (gray). c,f) surface cubes (gray boxes). d,e,g) several steps of floodfill propagation (starting at yellow circle, following yellow arrow). i) Corner points of each surface cube with exterior (yellow) or interior (red) state. j) Segments that cross the boundary of the output surfaces (Black lines). k) Intersection points (white circles) segments intersect the output surface. l) Lookup table of 3D surface constructions with different edge intersection patterns. m,n) Triangles (black lines) approximating output region (gray).

152 Next, on each cube of the output surface, we identify edges that connect one corner that is
 153 inside the output region to one that is outside (fig. 2j). Since these edges must pass through
 154 the output surface, we call *segIntersect()* on the root node to find the point of intersection
 155 between the edge and the output surface (fig. 2k). This process is parallelized across the list
 156 of edges, ensuring that the calculation is never duplicated when dealing with adjacent cubes.

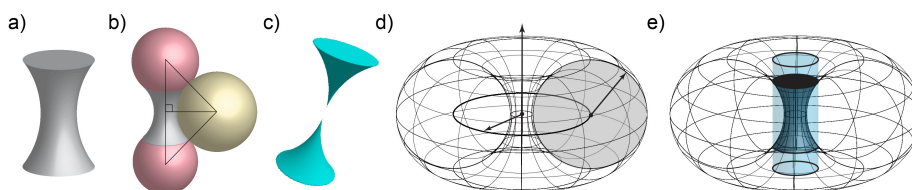
157 Finally, once intersections for every edge on every surface cube are determined, triangles
 158 are generated in each cube following a lookup table (fig. 2l). The collection of all resulting
 159 triangles form a closed triangular mesh that approximates the output region (fig. 2m,n).

2.2 Finding All Surface Cubes

160 `findAllSurfaceCubes()` accepts a cubic lattice (fig. 2b), a list of starting cubes (e.g. fig. 2c,f),
 161 and a CSG tree node for which to find all remaining surface cubes. We perform a parallel
 162 floodfill algorithm to find all remaining surface cubes: Each available thread is assigned a
 163 cube from the queue. Each thread tests cubes adjacent to the assigned cube to find any
 164 that are also on the surface of the input solid (fig. 2d). This test is performed by calling
 165 `containsPoint()` on the corners of the adjacent cube. If at least one corner is inside the input
 166 solid and another corner is outside, the adjacent cube is stored on a queue of upcoming
 167 cubes. Once all cubes adjacent to the initial surface cubes have been either added to the
 168 queue or discarded, all threads are then directed to find cubes adjacent to those still on the
 169 queue (e.g. fig. 2e), and so on, until the queue is empty, and all cubes on the surface of the
 170 input solid have been identified. Duplicate entries onto the queue are avoided by recording
 171 previously-examined cubes on a parallel hash table.
 172

2.3 Nodes of the CSG Tree

173 `pClay` supports several kinds of simple and complex solids for CSG operations. These are
 174 spheres, tetrahedra, spindles, and molecular surfaces. Our implementation of each type
 175 supports three basic functions: `containsPoint()`, `intersectSegment()`, and `findSurfaceCubes()`.
 176 To describe the implementation of these solids, we describe how each method is implemented
 177 for the solid. Spheres and tetrahedra are excluded because their implementations are trivial.
 178



■ **Figure 3** a) Spindle. b) Formation of a spindle (gray) from two atoms (red) and a solvent sphere (yellow). c) “Broken” spindle. d) Torus defining the characteristics of a spindle, including center point (black dot), perpendicular vector (vertical arrow), major radius (arrow from center point to ellipse), minor radius (arrow from ellipse to torus surface). e) Cylinder (light blue).

179 **Spindles** Spindles (fig. 3a) define the solvent excluded region between two atoms that are
 180 too close to permit a sphere representing a solvent molecule to pass between them (fig. 3b).
 181 “Broken” spindles (fig. 3c) can occur when the edge of the solvent sphere can pass beyond
 182 the centerline of the two atoms. Conceptually, spindles are the volume within a cylinder
 183 minus the volume within a coaxial torus. We define spindles by center point, perpendicular
 184 vector, major radius, and minor radius taken from the torus (fig. 3d), and end cap positions
 185 along the perpendicular vector (fig. 3e). The center point is the perpendicular projection
 186 of the center of the solvent sphere onto the segment between atom centers (fig. 3b). The
 187 perpendicular vector points from the center point towards the center of one atom. The major
 188 radius is the radius of the circle defined by the center of the solvent sphere as it rotates
 189 about the two atoms. The minor radius is the radius of the solvent sphere. The endcaps are
 190 circles perpendicular to the perpendicular vector that are defined by the point of tangency
 191 between the solvent sphere and the atoms, as the solvent sphere rotates about the atoms.
 192 The boundary surface of a spindle is defined by the end caps and elsewhere by the interior
 193 curve of the torus (fig. 3d).

194 To implement `containsPoint(p)`, note that the spindle is rotationally symmetric about the
 195 perpendicular vector. Thus, a plane K can be defined coplanar to p and the perpendicular
 196 vector of the torus. In K , p is inside the spindle only if it is inside the rectangle that
 197 defines the rotational cross section of the cylinder and also outside the circle that defines the
 198 rotational cross section of the torus.

199 intersectSegment(s) is computed by first setting up the calculation by translating the
 200 center of the spindle to the origin and rotating its axis to align it with the x axis. s is
 201 translated and rotated with it. We can describe the torus aligned to the x axis as

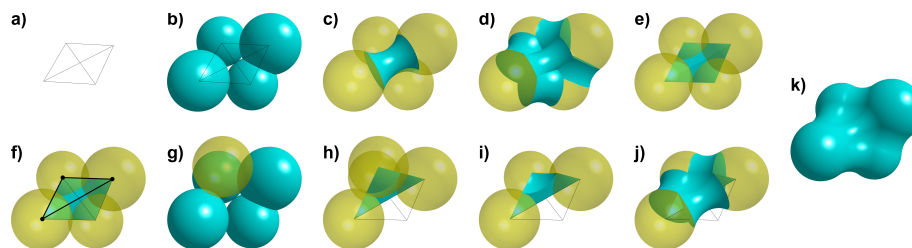
$$(x^2 + y^2 + z^2 + R^2 - r^2)^2 - 4R^2(y^2 + z^2) = 0$$

202 where R is the major radius, and r is the minor radius of the torus. In the torus equation,
 203 we substitute x , y and z with the line expressions $x_0 + td_x$, $y_0 + td_y$, and $z_0 + td_z$, where
 204 x_0, y_0, z_0 are segment starting points, and t parameterizes the line containing the line segment.
 205 The result of this substitution is a quartic equation on t , and roots of the equation will be
 206 parameters on the segment at points of intersection between the segment and the torus. We
 207 converted this equation into a monic quartic using Maxima, a computer algebra system [24].

208 To find the roots of this equation, we produce the Frobenius companion matrix of this
 209 quartic polynomial. The roots are the eigenvalues of this matrix. Here, complex eigenvalues
 210 will correspond to nonexistent points of intersection between the segment and the torus while
 211 real eigenvalues correspond to intersection points on the torus. We find these intersection
 212 points and eliminate any intersections that are outside of the cylinder. Separately, we also
 213 find intersections with the end caps of the spindle, treating them first as infinite planes and
 214 then determining if the intersection point is within the circle on the plane. Intersections
 215 between the segment and the endcaps or between the segment and the torus are returned as
 216 intervals where the segment is inside the spindle.

217 findStartingCubes() is implemented by first generating the segment between the centers
 218 of the endcaps. The lattice cube containing one centroid is identified, and if it is not a surface
 219 cube, the adjoining cube, through whose face which the segment passes, is identified as the
 220 next cube to examine. This process is repeated until either the segment ends at the other
 221 centroid of a surface cube has been found. In the case where the spindle is broken (fig. 3d),
 222 two segments are generated, starting at one endcap centroid and moving towards the other
 223 endcap centroid, but ending at the center.

224 **Molecular Solids** pClay generates molecular solids by positioning structural components
 225 with the power diagram [1]. This approach follows the classic methods for generating
 226 molecular surfaces, such as CASTp [31], MSMS [30], GRASP2 [29], which also use power
 227 diagrams or similar constructs. For this reason, we paraphrase our approach here, expanding
 228 on points that differ from classic methods. As in the earlier methods, our approach represents
 229 water molecules as solvent spheres, which can be of any given radius. By calling basic
 230 functions from simpler primitives, pClay achieves an efficient implementation of the basic
 231 functions for the entire molecular solid without describing it as a CSG operation of many
 232 individual primitives.



■ **Figure 4** Molecular Surface Construction. a) Dual graph of a power diagram on four atoms (black lines, points). b) Sphere primitives from atoms (teal). c) Atoms (yellow) with one spindle (teal). d) Atoms with all spindles from edges of the dual graph. e) Tetrahedron primitive (teal). f) One triangle of the dual graph (black lines, dots). g) Solvent sphere (yellow) tangent to three atoms. h) New tetrahedron (teal) with corners in the center of the three atoms of the triangle and the solvent sphere. i) Cup region inside the new tetrahedron (teal). j) Cup, shown with three adjacent spindles (teal) and three atoms of the triangle (yellow). k) Finished molecular solid.

233 We begin with an input file from the Protein Data Bank (PDB). Using atomic coordinates
234 and Van der Waals radii for each atom, we first compute a power diagram with REGTET [3].
235 The power diagram divides three dimensional space into cells corresponding to each atom of
236 the input. The size of a cell relates to the Van der Waals radius of the atom, through the
237 power function. Using the power diagram, we construct a topologically dual geometric graph
238 (fig. 4a), which has a vertex at the center of each atom and an edge between any vertices
239 that correspond to adjacent cells. This dual graph defines the location of the primitives that
240 will comprise the molecular solid. In sequential stages, we generate all primitives of the same
241 type in parallel, starting with sphere primitives, then spindles, tetrahedra, and so on.

242 At every vertex of the dual graph, we create sphere primitives with the appropriate Van
243 der Waals radius of each atom (fig. 4b). Next, we examine every edge on the dual graph and
244 generate a spindle between the atoms on at the endpoint of each edge, except for overlong
245 edges that are longer than the sum of Van set Waals radius of the endpoint atoms and
246 the diameter of the solvent sphere (fig. 4c,d). Once all spindles are completed, we identify
247 all tetrahedra in the dual graph that lack an overlong edge and we generate a tetrahedron
248 primitive for each one (fig. 4e).

249 Next, we identify triangles on the dual graph that are not between two tetrahedra (fig. 4f).
250 These triangles define triplets of atoms that may be on the molecular surface. To determine
251 whether the atoms are on the surface, we place a solvent sphere tangent to all three atoms
252 (fig. 4g). If the solvent sphere does not collide with any other atoms, we create a *negsphere*: a
253 sphere primitive in the tangent location in the same size as the solvent that describes a region
254 of the solvent outside the molecular surface. We also generate a tetrahedron with corners on
255 the triangle and at the center of the negsphere (fig. 4h). The region inside this tetrahedron
256 and outside the negsphere is both inside the solvent excluded region and not occupied by
257 spindles or atoms or other tetrahedra. We call this concave subset of a tetrahedron a *cup*
258 (fig. 4i), and describe cups as a negsphere-tetrahedron pair. The concave surface of the cup
259 is continuous with the three adjacent spindles and atoms (fig. 4j). Once all triangles that are
260 not between two tetrahedra have been examined for the presence of a cup, the combination
261 of spheres, spindles, tetrahedra and negspheres form a molecular solid (fig. 4k).

262 To support the three basic functions, we store all of these primitives in a data structure
263 for rapid range-based lookup. First, we generate a bounding box for each primitive. Next, we
264 generate a lattice of cubes, where each cube is 2 angstroms on a side. Finally, we associate
265 each primitive with all lattice cubes that intersect its bounding box. These associations
266 act as a hashing function that enables us to rapidly identify any primitives nearby a given
267 cube in the lattice. Since real molecules have finite atomic density, and since primitives are
268 constructed from atoms and between atoms, the number of primitives associated with any
269 cube is finite. As a result, a hashing function based on the lattice achieves algorithmically
270 constant time lookup of nearby primitives.

271 For `containsPoint(p)`, given a point p , if p is outside the coarse lattice, then we immediately
272 return false, because p must be outside the molecular surface. If not, we determine which
273 cube c of the coarse lattice contains p . Next, we identify all primitives associated with c . We
274 use the `containsPoint()` function of each associated primitive to determine if p is inside the
275 primitive. If p is inside a `negSphere`, then p is outside the molecular surface. if p is inside
276 any other primitives, then p is considered inside the molecular surface. If p is not inside any
277 primitives, it is outside.

278 For `intersectSegment(s)`, given a segment s , we generate a list of cubes C that contain
279 some interval of s . Next, we generate a list of primitives P associated with the cubes in C .
280 We then query each primitive p in the list P for an interval of intersection between p and s

281 using the `intersectSegment()` method of each primitive. The output intervals generated are
 282 the union of the intervals in tetrahedra, spindles and spheres minus the union of intervals
 283 inside negspheres.

284 For `findStartingCubes(l)`, during the construction of the molecular solid, we record the
 285 points of tangency between all negspheres and atom spheres. For each of these points, we
 286 identify the lattice cubes of l that contain them. We also generate starting cubes from all
 287 spindles and isolated spheres in the protein structure, calling `findStartingCubes()` on each of
 288 these primitives. From these cubes, we return only cubes that exhibit one corner inside and
 289 one outside the molecular solid.

290 **CSG Operations** CSG operation nodes are non-leaf nodes that represent the outcome of a
 291 CSG operation on its operand nodes. They fulfill the three basic functions by calling on its
 292 operand nodes, which we refer to as A and B in the text below.

293 Given a point p , the CSG union returns true only when `containsPoint(p)` returns true
 294 on at least one operand. The CSG intersection returns true only when `containsPoint(p)`
 295 returns true on both operands. The CSG difference between A and B returns true only when
 296 $A.containsPoint(p)$ is true and $B.containsPoint(p)$ is false.

297 For a given segment s , `intersectSegment(s)` on any CSG operation calls `intersectSegment(s)`
 298 on operands A and B , generating intervals a and b . When run on a CSG Union, Intersection
 299 or Difference, the output is, respectively, the union, intersection, or difference of a and b .

300 Given a cubic lattice l , calling `getSurfaceCubes(l)` on a CSG union, intersection, or
 301 difference returns the setwise union of cubes returned by calling $A.getSurfaceCubes()$ and
 302 $B.getSurfaceCubes()$. We always return a union of cubes because examining the union
 303 of cubes can avoid circumstances where a disconnected region in the final solid is lost.
 304 Performance profiling revealed that considering the union of all cubes is a minor cost in
 305 overall performance, except in artificially constructed cases that create many irrelevant cubes.

306 **Implementation Details** pClay is implemented in C and C++. Interprocess communication
 307 was built with Intel's Threading Building Blocks library. A C wrapper connects REGTET
 308 [3] to pClay. Benchmarks were run on a dual Xeon E5-2609 system with 8 cores at 2.5 Ghz
 309 and 32GB of ram and on a Xeon Phi 7120P with 61 cores at 1.2 Ghz and 16GB ram.

310 **Datasets Used** *Dataset A* is 100 nonredundant pdb structures from VAST [23], with BLAST
 311 p-value cutoff $10e-7$. *Dataset B* is 30 spheres, spindles and tetrahedra randomly generated in
 312 a cube with 10\AA sides. *Dataset C* is 14 binding cavities from a nonredundant subset of the
 313 trypsins (1a0j,1aks,1ane,1aq7,1bzx,1fn8,1hrw,1trn,2eek,2f91), chymotrypsins (1eq9,8gch) and
 314 elastases (1b0e,1elt). More info: www.cse.lehigh.edu/~chen/papers/WABI2019/appendix.pdf

315 **3 Experimental Results**

316 **3.1 Accuracy of molecular solid generation**

317 We compared the surfaces of molecular solids from pClay to surfaces made with the trollbase
 318 library, an established tool for molecular surface generation used in GRASP2 [29], VASP [8],
 319 VASP-E [5], and MarkUs [12]. pClay surfaces were created at 0.25\AA resolution to yield a
 320 similar number of triangles and thus a fairer comparison. Using proteins from dataset A,
 321 surfaces generated by pClay had an average of 197,718.54 points. From these points, we
 322 measured the distance to the closest point on the surface generated by trollbase for the same
 323 protein. On average, that distance was 0.00383\AA (std dev 0.0004\AA). This tiny deviation
 324 between pClay and trollbase surfaces demonstrates the accuracy of pClay surfaces.

3.2 Performance Comparison and Parallel Scaling

To our knowledge, VASP [5] is the only existing algorithm for comparing molecular solids using CSG. It lacks exact primitives or parallelism, but it can identify steric elements of protein structure that control specificity [8, 6, 13]. We compared the performance of pClay and VASP on the same CSG operations using Xeon (CPU) and Xeon Phi (PHI) processors.

Random Primitives First, we compared CSG performance on the union of primitives in dataset B, generating mesh outputs at resolutions 1.0Å, .5Å, .25Å and .125Å. Since VASP does not use primitives, it was provided triangle meshes of identical primitives. All CSG trees were balanced, but imbalanced trees had nearly identical runtimes (not shown for brevity).

On one CPU core, pClay required .113 seconds to compute the union at 1.0Å resolution. 9.492 seconds were required to compute the same union at .125Å resolution (fig. 5a). Increasing to 8 threads, runtime dropped to .03 seconds for unions at 1.0Å resolution, and 1.465 seconds to at .125Å. In contrast, single-threaded VASP required 3 seconds to compute the same union at 1.0Å, and 64 at .125Å. pClay far outperformed VASP on one thread.

On 8, 16, 32, and 60 PHI cores, which are slower than CPU cores, runtimes exhibited sublinear improvement (fig. 5b). Runtimes on coarser resolutions improved less than for finer resolutions. The difference in parallel speedup (fig. 5c) arises from small problem sizes at coarse resolutions, where communications and setup outweigh the advantages of parallelism.

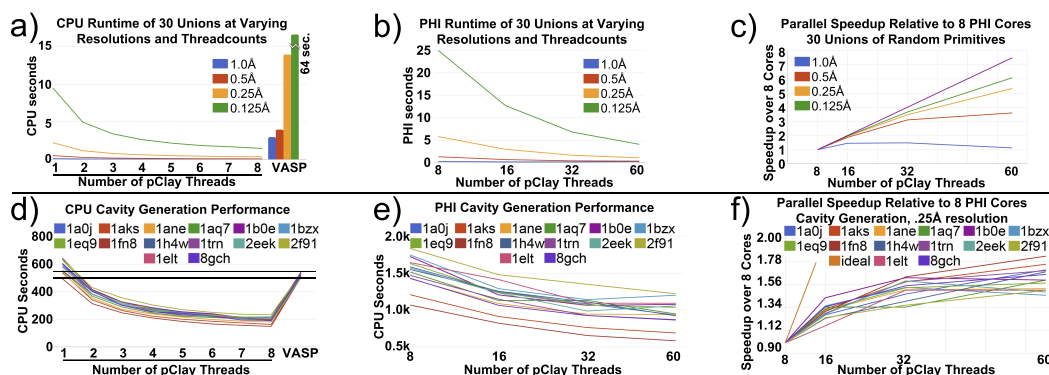


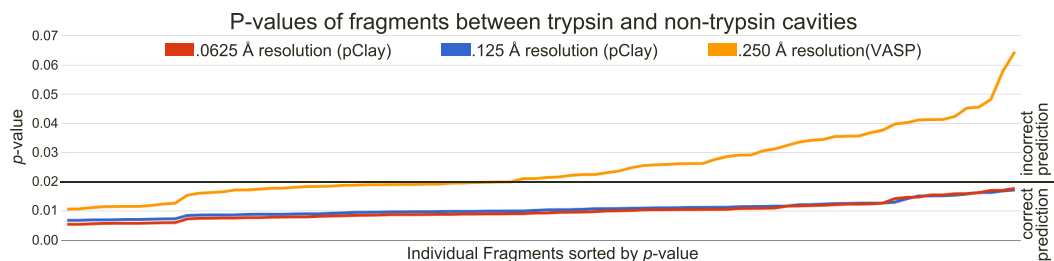
Figure 5 a) Time to compute the union of 30 random primitives at varying resolutions and CPU cores. VASP performance (single threaded) is shown in vertical bars. b) Time spent to compute the unions on PHI cores. c) Parallel speedup on PHI cores. d) Time spent for pClay to produce several binding cavities on CPU cores, compared to single-core VASP. e) Time to produce the same cavities on PHI cores. e) Parallel speedup of pClay in cavity production on varying PHI cores.

Binding Cavities We also tested pClay by generating binding cavities, using the method from fig. 1, on dataset C. While pClay is capable of much finer resolutions, all cavities were generated at .25Å, the practical resolution limit for VASP. Figure 5d plots cavity generation times for these cavities. pClay required between 493 and 643 seconds on one CPU core, and between 149 and 233 seconds on 8 cores. Single threaded VASP required between 499 and 538 seconds to perform the same work. Single threaded, pClay was slightly slower than VASP, but much faster when adding a second core, and faster still when adding more.

Cavity generation was also run on 8, 16, 32, and 60 Xeon Phi cores. Runtimes on PHI cores are slower than on CPU cores because PHI cores have slower clock speeds. Runtimes fell slowly as threads increased (fig. 5e). Substantial increases in the number of PHI cores resulted in only modest improvements in runtimes (fig. 5f). This result contrasted from those performed on CPU cores, where performance improved substantially with increases in parallelism. These results point to bottlenecks in the PHI architecture affected by cavity generation, which is more data intensive than unions of random primitives.

357 **3.3 Evaluating pClay on Existing Applications**

358 The added precision of pClay creates several new applications. We demonstrate one such
 359 application by producing training data for VASP-S, a statistical model for detecting differences
 360 in ligand binding specificity with steric causes [6]. VASP-S is trained on the volumes of
 361 individual CSG differences computed from cavities with the same binding preference. This
 362 training enables VASP-S to estimate the probability (the p -value) that two given cavities have
 363 similar binding preferences. If p is lower than a threshold α , VASP-S rejects the hypothesis
 364 that two cavities have similar binding preferences, and predicts that they are different.



■ **Figure 6** The p -value of the largest fragment between every trypsin-elastase and trypsin-chymotrypsin pair in Dataset C, estimated with training data generated at several resolutions (red, blue, orange lines). Fragments are sorted in ascending p -values along the horizontal axis. The black line indicates the α threshold of 0.02, below which we predict that the fragment is representative of proteins with different binding preferences. The finer-resolution training data, made possible with pClay, yielded more accurate predictions.

365 We hypothesized that training the VASP-S model with data generated at finer resolutions
 366 will produce more accurate predictions than a VASP-S model trained with coarser data.
 367 To test this hypothesis, we used cavities from the trypsins, which prefer to bind positively
 368 charged amino acids. These cavities contrast from those of the chymotrypsins and the
 369 elastases, which prefer large aromatics or small hydrophobics, respectively. Three training
 370 sets were constructed from these cavities by generating all possible CSG differences between
 371 all pairs of trypsins. VASP was used to produce a copy of the training set at 0.25 Å resolution
 372 and pClay was used to produce the same set at resolutions of 0.125 Å and 0.0625 Å. We then
 373 computed CSG differences between every trypsin and every nontrypsin at these resolutions.
 374 Finally, we estimate the p -value of the largest CSG difference between every trypsin and
 375 every non-trypsin in Dataset C, at all three resolutions (fig. 6). We expect VASP-S to
 376 produce a low p -value on these CSG differences.

377 Using the conservative α threshold of 2%, when trained at 0.25 Å resolution, VASP-S
 378 predicts that 43 of the 81 CSG differences between trypsin and non-trypsin cavities indicate
 379 binding preferences. This discrepancy indicates 38 false negative predictions, where VASP-S
 380 incorrectly overlooked cavities with different binding preferences. However, when trained
 381 at 0.125 Å or 0.0625 Å resolution, VASP-S correctly predicts that all CSG differences
 382 were from cavities with different binding preferences, a 0% false negative rate. These results
 383 demonstrate that pClay can provide superior precision, ensuring that existing aggregate
 384 methods do not lose accuracy by overlooking useful predictions.

385 **4 Discussion**

386 We have presented pClay, the first parallel algorithm for performing CSG analysis of protein
 387 structures at arbitrarily high resolutions, up to machine precision. It leverages mathematically
 388 exact primitives that can be assembled into molecular solids and parallel depth first search
 389 to compute CSG operations with multiple threads.

390 Molecular solids from pClay are nearly identical to molecular surfaces generated by
391 existing, widely used software. At hundreds of thousands of positions, pClay surfaces differed
392 from surfaces generated by existing methods by only thousandths of an angstrom. While
393 existing surface methods have generally been validated by visual examination, this exhaustive
394 comparison sets a new standard for validation. Surface generation stresses the algorithms
395 that underpin CSG operations, illustrating that pClay is making accurate comparisons.

396 We showed that pClay performs CSG operations efficiently on both artificial and realistic
397 data. Evaluating the method on both Xeon CPU and Xeon Phi architectures, pClay exhibited
398 scaleable multithreaded performance on all tests, though scaling was modest on the Xeon
399 Phi for cavity generation. These results show that parallelism can drive both efficiency and
400 precision for the comparison of protein structures.

401 Finally, we showed how the precision of pClay can advance existing methods by training
402 a statistical classifier to distinguish elements of protein structures that have a steric influence
403 on binding specificity. pClay provided training data to the classifier that was more precise
404 than what could have been provided by existing methods, enabling more accurate estimates
405 of statistical significance, and ultimately a total elimination of false negative predictions.

406 These capabilities enable applications in the detection and explanation of structural
407 features that influence binding preferences. For example, the statistical model tested here
408 finds elements of protein structures that could have a steric influence on specificity, thereby
409 generating an explanation based on a steric mechanism that typically requires human
410 expertise. As high throughput technologies increasingly reveal the ways in which disease
411 proteins can vary, pClay is a glimpse into a new space of techniques that can use protein
412 variants to supplement human experience in deciphering the structural mechanisms of
413 molecular recognition.

414 ——— References ———

- 415 1 F Aurenhammer. Power diagrams: properties, algorithms and applications. *SIAM Journal on*
416 *Computing*, 16(1):78–96, 1987.
- 417 2 CL Bajaj, V Pascucci, A Shamir, RJ Holt, and AN Netravali. Dynamic maintenance and
418 visualization of molecular surfaces. *Discrete Applied Mathematics*, 127(1):23–51, 2003.
- 419 3 J Bernal. REGTET: A program for computing regular tetrahedralizations. In *International*
420 *Conference on Computational Science*, pages 629–632. Springer, 2001.
- 421 4 M Brylinski and J Skolnick. A threading-based method (findsite) for ligand-binding site
422 prediction and functional annotation. *P Natl Acad Sci USA*, 105(1):129–134, 2008.
- 423 5 BY Chen. VASP-E: Specificity annotation with a volumetric analysis of electrostatic isopotentials.
424 *PLoS computational biology*, 10(8):e1003792, 2014.
- 425 6 BY Chen and S Bandyopadhyay. VASP-S: A volumetric analysis and statistical model
426 for predicting steric influences on protein-ligand binding specificity. In *IEEE Int Conf*
427 *Bioinformatics Biomed 2011*, pages 22–29. IEEE, 2011.
- 428 7 BY Chen, VY Fofanov, DH Bryant, BD Dodson, DM Kristensen, AM Lisewski, M Kimmel,
429 O Lichtarge, and LE Kavraki. The MASH pipeline for protein function prediction and an
430 algorithm for the geometric refinement of 3d motifs. *J Comput Biol*, 14(6):791–816, 2007.
- 431 8 BY Chen and B Honig. VASP: a volumetric analysis of surface properties yields insights into
432 protein-ligand binding specificity. *PLoS computational biology*, 6(8):e1000881, 2010.
- 433 9 ML Connolly. Analytical molecular surface calculation. *Journal of applied crystallography*,
434 16(5):548–558, 1983.
- 435 10 L Ellingson and J Zhang. Protein surface matching by combining local and global geometric
436 information. *PloS one*, 7(7):e40540, 2012.
- 437 11 F Ferre, G Ausiello, A Zanzoni, and M Helmer-Citterich. Functional annotation by identification
438 of local surface similarities: a novel tool for structural genomics. *BMC Bioinf*, 6(1):194, 2005.

- 439 12 M Fischer, QC Zhang, F Dey, BY Chen, B Honig, and D Petrey. Markus: a server to navigate
440 sequence–structure–function space. *Nucleic acids research*, 39(suppl_2):W357–W361, 2011.
- 441 13 BG Godshall and BY Chen. Improving accuracy in binding site comparison with homology
442 modeling. In *IEEE Int Conf Bioinformatics Biomed 2012*, pages 662–669. IEEE, 2012.
- 443 14 L He, F Vandin, G Pandurangan, and C Bailey-Kellogg. Ballast: a ball-based algorithm for
444 structural motifs. *Journal of Computational Biology*, 20(2):137–151, 2013.
- 445 15 C-L Hung and Y-L Lin. Implementation of a parallel protein structure alignment service on
446 cloud. *International journal of genomics*, 2013, 2013.
- 447 16 T Ju, F Losasso, S Schaefer, and J Warren. Dual contouring of hermite data. In *ACM*
448 *transactions on graphics (TOG)*, volume 21, pages 339–346. ACM, 2002.
- 449 17 F Kaiser, A Eisold, and D Labudde. A novel algorithm for enhanced structural motif matching
450 in proteins. *Journal of Computational Biology*, 22(7):698–713, 2015.
- 451 18 K Kinoshita, Y Murakami, and H Nakamura. eF-seek: prediction of the functional sites of
452 proteins by searching for similar electrostatic potential and molecular surface shape. *Nucleic*
453 *acids research*, 35(suppl_2):W398–W402, 2007.
- 454 19 J Konc, M Depolli, R Trobec, K Rozman, and D Janežič. Parallel-probis: Fast parallel
455 algorithm for local structural comparison of protein structures and binding sites. *Journal of*
456 *computational chemistry*, 33(27):2199–2203, 2012.
- 457 20 B Lee and FM Richards. The interpretation of protein structures: estimation of static
458 accessibility. *Journal of molecular biology*, 55(3):379–IN4, 1971.
- 459 21 Y Liu, K Shah, F Yang, L Witucki, and KM Shokat. A molecular gate which controls unnatural
460 atp analogue recognition by the tyrosine kinase v-src. *Bioorg Med Chem*, 6(8):1219–1226, 1998.
- 461 22 WE Lorensen and HE Cline. Marching cubes: A high resolution 3d surface construction
462 algorithm. In *ACM siggraph computer graphics*, volume 21, pages 163–169. ACM, 1987.
- 463 23 Thomas Madej, Christopher J Lanczycki, Dachuan Zhang, Paul A Thiessen, RC Geer,
464 A Marchler-Bauer, and SH Bryant. Mmdb and vast+: tracking structural similarities between
465 macromolecular complexes. *Nucleic acids research*, 42(D1):D297–D303, 2013.
- 466 24 WA Martin and RJ Fateman. The mcsyma system. In *Proceedings of the second ACM*
467 *symposium on Symbolic and algebraic manipulation*, pages 59–75. ACM, 1971.
- 468 25 NL Max and ED Getzoff. Spherical harmonic molecular surfaces. *IEEE Computer Graphics*
469 *and Applications*, 8(4):42–50, 1988.
- 470 26 EC Meng, BJ Polacco, and PC Babbitt. 3d motifs. In *From Protein Structure to Function*
471 *with Bioinformatics*, pages 187–216. Springer, 2009.
- 472 27 RJ Morris, RJ Najmanovich, A Kahraman, and JM Thornton. Real spherical harmonic
473 expansion coefficients as 3d shape descriptors for protein binding pocket and ligand comparisons.
474 *Bioinformatics*, 21(10):2347–2355, 2005.
- 475 28 BE Nolan, E Levenson, and BY Chen. Influential mutations in the smad4 trimer complex
476 can be detected from disruptions of electrostatic complementarity. *Journal of Computational*
477 *Biology*, 24(1):68–78, 2017.
- 478 29 D Petrey and B Honig. GRASP2: visualization, surface properties, and electrostatics of
479 macromolecular structures and sequences. *Methods in enzymology*, 374:492–509, 2003.
- 480 30 MF Sanner, AJ Olson, and J-C Spehner. Reduced surface: an efficient way to compute
481 molecular surfaces. *Biopolymers*, 38(3):305–320, 1996.
- 482 31 W Tian, C Chen, and J Liang. Castp 3.0: Computed atlas of surface topography of proteins
483 and beyond. *Biophysical Journal*, 114(3):50a, 2018.
- 484 32 J Venkateswaran, B Song, T Kahveci, and C Jermaine. Trial: A tool for finding distant
485 structural similarities. *IEEE/ACM Trans Comput Biol Bioinform*, 8(3):819–831, 2011.
- 486 33 Y Zhou, X-P Li, BY Chen, and NE Tumer. Ricin uses arginine 235 as an anchor residue to
487 bind to p-proteins of the ribosomal stalk. *Scientific reports*, 7:42912, 2017.