

1 **Title:**

2 Predicting Allosteric Hotspots Using Dynamics-Based Formalisms with Sequence
3 Analyses Across Diverse Evolutionary Timescales

4
5 **Authors & associated information:**

6 Declan Clarke^{a,1}, Anurag Sethi^{b,c,1}, Shantao Li^{b,d}, Sushant Kumar^{b,c}, Richard W.F.
7 Chang^e, Jieming Chen^{b,f}, and Mark Gerstein^{b,c,d,2}

8
9 ^aDepartment of Chemistry, Yale University, 260/266 Whitney Avenue PO Box 208114,
10 New Haven, CT 06520 USA

11 ^bProgram in Computational Biology and Bioinformatics, Yale University, 260/266
12 Whitney Avenue PO Box 208114, New Haven, CT 06520, USA

13 ^cDepartment of Molecular Biophysics and Biochemistry, Yale University, 260/266
14 Whitney Avenue PO Box 208114, New Haven, CT 06520, USA

15 ^dDepartment of Computer Science, Yale University, 260/266 Whitney Avenue PO Box
16 208114, New Haven, CT 06520, USA

17 ^eYale College, 260/266 Whitney Avenue PO Box 208114, New Haven, CT 06520, USA

18 ^fIntegrated Graduate Program in Physical and Engineering Biology, Yale University,
19 260/266 Whitney Avenue PO Box 208114, New Haven, CT 06520, USA

20
21 ¹ D.C. and A.S. contributed equally to this work.

22 ² Correspondence should be addressed to M.G. (pi@gersteinlab.org)

DECLAN CLARKE 12/8/15 3:18 PM

Deleted:

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: Identifying allosteric hotspots with dynamics: application to conservation in deep sequencing -

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Font:Not Bold

ABSTRACT

43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64

The rapidly growing volume of data being produced by next-generation sequencing initiatives is enabling more in-depth analyses of conservation than previously possible. Deep sequencing is uncovering disease loci and regions under selective constraint, despite the fact that intuitive biophysical reasons for such constraint are sometimes unavailable. Allostery may often provide the missing explanatory link. We use models of protein conformational change to identify allosteric residues, by predicting essential surface cavities or information flow bottlenecks, and we develop a software tool (stress.molmovdb.org) that enables users to perform this analysis on their own proteins of interest. Though fundamentally 3D-structural in nature, this software is computationally fast, thereby allowing us to run it across the PDB and to evaluate general properties of predicted allosteric residues, which tend to be conserved over long and short evolutionary time scales. We highlight examples in which allosteric residues can help explain poorly understood disease-associated variants.

MISSING

FINDING

OUR ANALYSIS

OF THAT

WE FIND THESE

- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: protein
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: protein
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: strong
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: , in many cases, we cannot find
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: (such as the need to engage in protein-protein interactions or to achieve a close-packed hydrophobic core). Allosteric hotspots
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: Here, we
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: such
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: In particular, we predict allosteric residues that can act as
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: While our tool is
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: it
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: still
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: and tractable. This allows
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: entire Protein Data Bank
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: large-scale
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: . We find that these
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: significantly
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: both
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: Finally, we
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: specific
- DECLAN CLARKE 12/8/15 3:18 PM
Deleted: previously

89 INTRODUCTION

90 The ability to sequence large numbers of human genomes is providing a much
91 deeper view into protein evolution. When trying to understand the evolutionary pressures
92 on a given protein, structural biologists now have at their disposal an unprecedented
93 breadth of data regarding patterns of conservation, both across species and amongst
94 humans. As such, there are greater opportunities to take a ~~more~~^{more} integrated view of the
95 context in which a protein and its residues function. This ~~integrated~~ view necessarily
96 includes structural constraints such as residue packing, protein-protein interactions, and
97 stability. However, deep sequencing is unearthing a class of conserved residues on which
98 no obvious structural constraints appear to be acting. The missing link in understanding
99 these regions may ~~can~~ be provided by considering the protein's dynamic behavior and
100 distinct functional states within an ensemble.

101 The underlying energetic landscape responsible for the relative distributions of
102 alternative conformations is dynamic in nature: allosteric signals or other external
103 changes may reconfigure and reshape the landscape, thereby shifting the relative
104 populations of states within an ensemble (Tsai et al., 1999). Landscape theory thus
105 provides the conceptual underpinnings necessary to describe how proteins change
106 behavior and shape under changing conditions. A primary driving force behind the
107 evolution of these landscapes is the need to efficiently regulate activity in response to
108 changing cellular contexts, thereby making allostery and conformational change essential
109 components of protein evolution.

110 Given the importance of allosteric regulation, as well as ~~the role of allostery~~ ^{its} in
111 imparting efficient functionality, several methods have been devised for the identification
112 of likely allosteric residues. Conservation itself has been used, either in the context of
113 conserved residues (Panjkovich and Daura, 2012), networks of co-evolving residues
114 (Halabi et al., 2009; Lee et al., 2008; Lockless et al., 1999; Reynolds et al., 2011;
115 Shulman et al., 2004; Süel et al., 2003), or local conservation in structure (Panjkovich
116 and Daura, 2010). In related studies, both conservation and geometric-based searches for
117 allosteric sites have been successfully applied to several systems (Capra et al., 2009). A
118 number of methods employing support vector machines ^{RSL} have also been described (Huang
119 and Schroeder, 2006; Huang et al., 2013). Normal modes analysis, coupled with ligands
120 of varying size, have been used to examine the extent to which bound ligands interfere
121 with low-frequency motions, thereby identifying potentially important residues at the
122 surface (Ming and Wall, 2005; Mitternacht and Berezovsky, 2011; Panjkovich and
123 Daura, 2012).

124 The concept of 'protein quakes' has been introduced to explain local regions of
125 proteins that are essential for conformation transitions (Miyashita et al., 2003). A protein
126 may relieve the strain of a high-energy configuration by local structural changes. Such
127 local changes often occur at the focal points of allosteric regulation, and these regions
128 may be identified in a number of ways, including modified normal modes analysis
129 (Miyashita et al., 2003) or time-resolved X-ray scattering (Arnlund et al., 2014).

130 Normal modes have also been used by the Bahar group to identify important
131 subunits that act in a coherent manner for specific proteins (Chennubhotla and Bahar,
132 2006; Yang and Bahar, 2005). Rodgers *et al* have applied normal modes to identify key

IN MORE
3D...

133 residues in CRP/FNR transcription factors (Rodgers et al., 2013). Molecular dynamics
134 (MD) and network analyses have been used to identify interior residues that may function
135 as allosteric bottlenecks (Csermely et al., 2013; Gasper et al., 2012; Rousseau and
136 Schymkowitz, 2005; Sethi et al., 2009; Vanwart et al., 2012). Along similar lines, Ghosh
137 et al. (2008) have taken a novel approach of combining MD and network principles to
138 characterize allosterically important inter-domain communication in methionyl tRNA
139 synthetase (Munro, 2009). In conjunction with NMR, Rivalta *et al* use MD and
140 network analysis to identify important regions in imidazole glycerol phosphate synthase
141 (Rivalta et al., 2012).

142 Though having provided valuable insights, many of these approaches may be
143 limited in terms of scale (the numbers of proteins which may feasibly be investigated),
144 computational demands, or the class of residues to which the method is tailored (surface
145 or interior). Using models of protein conformational change, ^{here} we identify both surface and
146 interior residues that may act as essential allosteric regions in a computationally tractable
147 manner, thereby enabling high-throughput analysis. This framework directly incorporates
148 information regarding ^{3D} protein structure and dynamics, and it is ^{CHANGE} applied to proteins
149 throughout the PDB (Berman et al., 2000) that exhibit conformational change. The
150 relatively greater conservation of the residues identified (both across species and amongst
151 humans) may help to elucidate many of the otherwise poorly understood regions in
152 proteins. In a similar vein, several of our identified sites correspond to human disease loci
153 for which no clear mechanism for pathogenesis had previously been proposed. Finally, ^{WE MAKES PLURAL.}
154 our framework (termed STRESS, for STRucturally-identified ESSential residues) is made
155 available through a tool to enable users to submit their own structures for analysis.

156

157

RESULTS

158

Identifying Potential Allosteric Residues

159

Allosteric residues at the surface generally play a regulatory role that is

160

fundamentally distinct from that of allosteric residues within the protein interior. While

161

surface residues may often constitute the sources or sinks of allosteric signals, interior

162

residues act to transmit such signals. We use models of protein conformational change in

163

an attempt to identify both classes of residues (Figure 1). Throughout, we term these

164

potential allosteric residues at the surface and interior “surface-critical” and “interior-

165

critical” residues, respectively. Critical residues are first identified in a set of 12 well-

166

studied canonical systems (see Figure S1, as well as Table S1 for rationale), and they are

167

then identified on a large scale across hundreds of distinct proteins.

168

Identifying Surface-Critical Residues

169

Allosteric ligands often act by binding to surface cavities and modulating protein

170

conformational dynamics. The surface-critical residues, some of which may act as latent

171

ligand binding sites and active sites, are first identified by finding cavities using Monte

172

Carlo simulations to probe the surface with a flexible ligand (Figure 1A, top-left). The

173

degree to which cavity occlusion by ~~the~~^a ligand disrupts large-scale conformational

174

change is used to assign a score to each cavity – sites at which ligand occlusion strongly

175

interfere with conformational change earn high scores (Figure 1A, top-right), whereas

176

shallow pockets (Figure 1A, bottom-left) or sites at which large-scale motions are largely

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Indent: First line: 0"

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: ... [1]

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: for which both the *holo* and *apo* states are available (Table S1 and

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: .

182 unaffected (Figure 1A, bottom-right) earn lower scores. Further details are [provided](#) in SI
183 Methods [section 3.1-a](#).

184 This approach is a modified version of the binding leverage framework
185 introduced by Mitternacht and Berezovsky (Mitternacht and Berezovsky, 2011). The
186 main modifications include the use of heavy atoms in the protein during the Monte Carlo
187 search, in addition to an automated means of thresholding the list of ranked scores. These
188 modifications were implemented to provide a more selective set of sites. Without them,
189 an exceedingly large fraction of the protein surface would be captured (Figure 2C). We
190 find that this modified approach results in an average of ~2 distinct sites per domain
191 (Figure 2A). The distribution for distinct sites within entire complexes is given in Figure
192 2B.

193 Within the canonical set of 12 proteins, we positively identify an average of 56%
194 of the sites known to be directly involved in ligand or substrate binding (see [Table 1](#),
195 Figure S1, and [SI Methods section 3.1-a-iv](#)). Some of the sites identified do not directly
196 overlap with known binding regions, but we often find that these “false positives”
197 nevertheless exhibit some degree of overlap with binding sites ([Table S2](#)). In addition,
198 those surface-critical sites that do not match known binding sites may nevertheless
199 correspond to latent allosteric regions: even if no known biological function is assigned
200 to such regions, their occlusion may nevertheless disrupt large-scale motions.

HITHERTO UNFOUND

202 **Dynamical Network Analysis to Identify Interior-Critical Residues**

203 The binding leverage framework described above is intended to capture hotspot
204 regions at the protein surface, but the Monte Carlo search employed is *a priori* excluded
205 from the protein interior. Allosteric residues often act within the protein interior by

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: (see SI Methods).

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: (see SI Methods).

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S2

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: ; see SI Methods for details on defining distinct sites).

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 60

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: Tables S2 and S3

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: supplementary note “Capturing Known Ligand-Binding Sites”).

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S4

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Font:Arial, 11 pt, Bold

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: -

... [2]

218 functioning as essential ‘bottlenecks’ within the communication pathways between distal
219 regions. An allosteric signal transmitted from one region to another may conceivably take
220 various alternative routes, but many of these routes can share a common set of residues.
221 The removal of such a common set of residues can result in the loss of many or all of the
222 available routes for allosteric signal transmission, thereby making these residues essential
223 information flow bottlenecks.

224 To identify bottlenecks, the protein is first modeled as a network, wherein
225 residues represent nodes and edges represent contacts between residues (in much the
226 same way that the protein is modeled as a network in constructing anisotropic network
227 models, see below). In this regard, the problem of identifying interior-critical residues is
228 reduced to a problem of identifying nodes that participate in network bottlenecks (see
229 Figure 1B and SI Methods [section 3.1-b](#) for details). Briefly, the network edges are first
230 weighted by the ^{DEG OF V STREN. OF...} correlated motions of contacting residues: a strong correlation in the
231 motion between contacting residues implies that knowing how one residue moves better
232 enables one to predict the motion of the other, thereby suggesting a strong information
233 flow between the two residues. The weights are used to assign ‘effective distances’
234 between connecting nodes, with strong correlations resulting in shorter effective node-
235 node distances.

236 Using the motion-weighted network, “communities” of nodes are identified using
237 the Girvan-Newman formalism (Girvan et al., 2002). A community is a group of nodes
238 such that each node within the community is highly inter-connected, but loosely
239 connected to other nodes outside the community. Communities are thus densely inter-

240 connected regions within proteins. As tangible examples, the community partitions and
241 the resultant critical residues for the canonical set are given in Figures S2.

242 Finally, the betweenness of each edge is calculated. The betweenness of an edge
243 is defined as the number of shortest paths between all pairs of residues that pass through
244 that edge, with each path representing the sum of effective node-node distances assigned
245 in the weighting scheme above. Those residues that are involved in the highest-
246 betweenness edges between pairs of interacting communities are identified as the
247 interior-critical residues. These residues are essential for information flow between
248 communities, as their removal would result in substantially longer paths between the
249 residues of one community to those of another.

250

251 **Software Tool: STRESS (STRucturally-identified ESSential residues)**

252 The implementations for finding both surface- and interior-critical residues have
253 been made available to the scientific community through a new software tool, STRESS,
254 which may be accessed at stress.molmovdb.org (Figure 3A). Users may specify a PDB to
255 be analyzed, and the output provided constitutes the set of identified critical residues.

256 Obviating the need for long wait times, the algorithmic implementation of our
257 software is highly efficient (Figures 3B and 3C). A typical protein of ~500 residues takes
258 only about 30 minutes on a 2.6GHz CPU. Running times are also minimized by using a
259 scalable server architecture that runs on the Amazon cloud (Figure 3D). A light front-end
260 server handles incoming user requests, and more powerful back-end servers, which
261 perform the calculations, are automatically and dynamically scalable, thereby ensuring
262 that they can handle varying levels of demand both efficiently and economically.

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S3 and S4

PASSIVE
WE MAKE

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S5

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: Figure S6

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S7

OVERALL

267 **High-Throughput Identification of Alternative**
268 **Conformations**

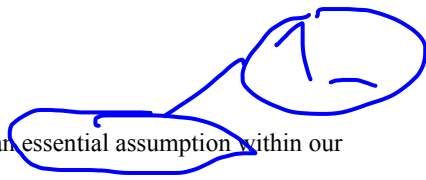
269 Pronounced conformational change is an essential assumption within our
270 framework for identifying potential allosteric residues. We use a generalized approach to
271 systematically identify instances of alternative conformations within the PDB. We first
272 perform multiple structure alignments (MSAs) across sequence-identical proteins that are
273 pre-filtered to ensure structural quality. We then use the resultant pairwise RMSD values
274 to infer distinct conformational states (Figure S3; see also SI Methods section 3.2).

275 The distributions of the resultant numbers of conformations for domains and
276 chains are given in Figures S3D and S3E, respectively, and an overview is given in
277 Figure S3F. We note that the alternative conformations identified arise in an extremely
278 diverse set of biological contexts, including conformational transitions that accompany
279 ligand binding, protein-protein or protein-nucleic acid interactions, post-translational
280 modifications, changes in oxidation or oligomerization states, etc. The dataset of
281 alternative conformations identified is provided as a resource in File S1 (see also Figure
282 S3G).

284 **Evaluating Conservation of Critical Residues**
285 **Using Various Metrics and Sources of Data**

286 The large number of dynamic proteins culled throughout the PDB, coupled with
287 the high algorithmic efficiency of our critical residue search implementation, provide a

TRANS



DECLAN CLARKE 12/8/15 3:18 PM
Deleted: Figures S8 and S9

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: for details

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: S2C

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: S2D

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: 2E. Further summary statistics are provided in Figure S10

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: (Figure S11).

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: S12

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: ... [3]

DECLAN CLARKE 12/8/15 3:18 PM
Formatted: Indent: First line: 0.5"

DECLAN CLARKE 12/8/15 3:18 PM
Formatted: Font:Times New Roman, No underline

LARGE pool...

pool,

298 means of evaluating general properties of these residues on a large scale. In particular, we
299 measure their conservation, as evaluated both over long (inter-species) and short (intra-
300 human) evolutionary timescales. Using a variety of conservation metrics and sources of
301 data, we find that both surface-critical (Figures 4A-D) and interior-critical (Figures 4E-H)
302 are consistently more conserved than non-critical residues. We emphasize that the
303 signatures of conservation identified not only provide a means of rationalizing many of
304 the otherwise poorly understood regions of proteins, but they also reinforce the functional
305 importance of the residues believed to be allosteric.

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: 3A
DECLAN CLARKE 12/8/15 3:18 PM
Deleted: 3E

307 Conservation Across Species

308 When evaluating conservation across species, we find that both surface- and
309 interior-critical residues tend to be significantly more conserved than non-critical residues
310 with the same degree of burial (Figures 4B and 4F, respectively). Surface-critical residue
311 sets have a mean conservation score (i.e., ConSurf score, see SI Methods section 3.3-a) of
312 -0.131, whereas non-critical residue sets with the same degree of burial have a mean
313 score of +0.059 ($p < 2.2e-16$; negative conservation scores designate stronger
314 conservation). Interior-critical residues exhibit a similar trend: the mean conservation
315 score for interior-critical residues and non-critical residues with the same degree of burial
316 is -0.179 and -0.102, respectively ($p=3.67e-11$).

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: 3B
DECLAN CLARKE 12/8/15 3:18 PM
Deleted: 3F

CONDENSE
CAPTION

317 Measures of Conservation Amongst Humans from Next-Generation Sequencing

318 We may also use sequenced human genomes and exomes to investigate
319 conservation, as many constraints may be human-specific and active in more recent
320 evolutionary history. In this context, commonly used metrics for evaluating conservation
321 include minor allele frequency (MAF) and derived allele frequency (DAF). Low MAF or

326 DAF values are interpreted as signatures of deleteriousness, as purifying selection is
327 prone to reduce the frequencies of harmful variants (see SI Methods [section 3.3-b](#)).

328 We find that 1000 Genomes (McVean et al., 2012) single-nucleotide variants
329 (SNVs) that hit surface-critical residues tend to occur at lower DAF values (Figure [4C](#)).

330 Though not significant, the significance improves when examining the shift in DAF
331 distributions, as evaluated with a KS test ($p=0.159$, Figure [S4A](#)), and we point out the
332 limited number of proteins (thirty-two) in which 1000 Genomes SNVs hit these critical
333 sites. Furthermore, the long tail extending to lower DAF values for surface-critical
334 residues may suggest that only a subset of the residues in our prioritized binding sites is
335 essential. However 1000 Genomes SNVs tend to hit interior-critical residues with
336 significantly lower DAF values than non-critical residues (Figure [4G](#); see also Figure
337 [S4B](#)).

338 Given the relatively small number of proteins to be hit by 1000 Genomes SNVs,
339 we also analyzed data provided by the Exome Aggregation Consortium (ExAC,
340 Cambridge MA 2015). ExAC provides sequence data for many more individuals, and the
341 ExAC sequencing itself is performed at much higher coverage. Thus, using MAF as a
342 conservation metric, we performed a similar analysis using this data. MAF distributions
343 for surface- and non-critical residues in the same set of proteins are given in Figure [4D](#).
344 Although the mean value of the MAF distribution for surface-critical residues is slightly
345 higher than that of non-critical residues, the median for surface-critical residues is
346 substantially lower than that for non-critical residues, demonstrating that the majority of
347 proteins are such that MAF values are lower in surface- than in non-critical residues. In

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 3C

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 13A

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: (see SI Methods).

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 3G

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S13B

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 3D

354 addition, the overall shifts of these distributions also point to a trend of lower MAF
355 values in surface-critical residues (Figure [S5A](#), KS test $p=9.49e-2$).

356 Interior-critical residues exhibit significantly lower MAF values than do non-
357 critical residues in the same set of proteins. MAF distributions for interior- and non-
358 critical residues are given in Figure [4H](#) (see also Figure [S5B](#)).

359 In addition to overall allele frequency distributions, one may also evaluate the
360 *fraction* of rare alleles as a metric for measuring selective pressure. This fraction is
361 defined as the ratio of the number of low-DAF or low-MAF (i.e., rare) non-synonymous
362 SNVs to all non-synonymous SNVs in a given protein annotation (such as all surface-
363 critical residues of the protein, for example; see SI Methods [section 3.3-b](#)). A higher
364 fraction is interpreted as a proxy for greater conservation (Khurana et al., 2013). Using
365 variable DAF (MAF) cutoffs to define rarity for 1000 Genomes (ExAC) SNVs, both
366 surface- and interior-critical residues are shown to harbor a higher fraction of rare alleles
367 than do non-critical residues, further suggesting a greater degree of evolutionary
368 constraint in critical residues (See [Figure 5](#)).

370 Comparisons Between Different Models of Protein Motions

371 Conformational changes may be modeled using vectors connecting pairs of
372 corresponding residues in crystal structures from alternative conformations (we term this
373 approach “ACT”, for “absolute conformational transitions”; see SI Methods [section 3.2-](#)
374 [c](#)). The crystal structures of such paired conformations may be obtained using the
375 framework discussed above. The protein motions may also be inferred from anisotropic
376 network models (ANMs) (Atilgan et al., 2001). ANMs entail modeling interacting
377 residues as nodes linked by flexible springs, in a manner similar to elastic network

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: S14A

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: 3H

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: S14B

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: Figure S15 and S16 for 1000
Genomes and ExAC data, respectively

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: ”).

DECLAN CLARKE 12/8/15 3:18 PM
Deleted: (further details in SI Methods).

REV.

TRANS
SUMMARY
"WE
GET
VECT IN
MANY
WAYS..."

385 models (Fuglebakk et al., 2015; Tirion, 1996) or normal modes analysis (Figure 1B).
386 ANMs are not only simple and straightforward to apply on a database scale, but unlike
387 using alternative crystal structures, the motion vectors inferred may be generated using a
388 single structure, and we thus use ANMs as our primary means of inferring motions.

389 Using vectors from either ACTs or ANMs give the same general results in terms
390 of the disparities in conservation between critical and non-critical residues. This method
391 is thus general with respect to how motion vectors are defined (see Figure 6 and SI
392 [Methods section 3.2-c](#) for further details).

394 **Critical Residues in the Context of Human Disease Variants**

395 Directly related to conservation is the concept of SNV deleteriousness: changes in
396 amino acid composition at specific loci may be more or less likely to result in disease.
397 SIFT (Ng and Henikoff, 2001) and PolyPhen (Adzhubei et al., 2010) are two tools for
398 predicting such effects, and we evaluated these predictions for critical and non-critical
399 residues hit by SNVs in ExAC. SNVs hitting critical residues exhibit significantly higher
400 PolyPhen scores relative to non-critical residues, suggesting the potentially higher disease
401 susceptibility at critical residues (Figure S6), though such significant disparities were not
402 observed in SIFT scores (Figure S7).

403 Using HGMD (Stenson et al., 2014) and ClinVar (Landrum et al., 2014), we
404 identify proteins with critical residues that coincide with disease-associated SNVs (Figure
405 [7A](#) and File S2). Several critical residues coincide with known disease loci for which the
406 mechanism of pathogenicity is otherwise unclear (File S3). The fibroblast growth factor
407 receptor (FGFR) is a case-in-point (Figure 7). SNVs in FGFR have been linked to
408 craniofacial defects. Dotted lines in Figure 7B highlight poorly understood disease SNVs

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S17

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: Supplemental note "Modeling Protein Motions by Directly Using Displacement Vectors from Alternative Conformations"

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S18

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S19

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 4A

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 4

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 4B

DISC
WORDS

419 that coincide with critical residues. In addition, we identify Y328 as a surface-critical
420 residue, which coincides with a disease-associated SNV from HGDM, despite the lack of
421 confident predictions of deleteriousness by several widely used tools for predicting
422 disease-associated SNVs, including PolyPhen (Adzhubei et al., 2010), SIFT (Ng and
423 Henikoff, 2001), and SNPs&GO (Calabrese et al., 2009). Together, these results suggest
424 that the incorporation of surface- and interior-critical residues introduces a valuable layer
425 of annotation to the protein sequence, and may help to explain otherwise poorly
426 understood disease-associated SNVs.

427

428 **DISCUSSION & CONCLUSIONS**

429 The same principles of energy landscape theory that dictate protein folding are
430 integral to how proteins explore different conformations once they adopt their folded
431 states. These landscapes are shaped not only by the protein sequence itself, but also by
432 extrinsic conditions. Such external factors often regulate protein activity by introducing
433 allosteric-induced changes, which ultimately reflect changes in the shapes and population
434 distributions of the energetic landscape. In this regard, allostery provides an ideal
435 platform from which to study protein behavior in the context of their energetic
436 landscapes. To investigate allosteric regulation, and to simultaneously add an extra layer
437 of annotation in the context of conservation patterns, an integrated framework to identify
438 potential allosteric residues is essential. We introduce a framework to select such
439 residues, using knowledge of conformational change.

440 When applied to many proteins with distinct conformational changes in the PDB,
441 we investigate the conservation of potential allosteric residues in both inter-species and
442 intra-human genomes contexts, and find that these residues tend to exhibit greater
443 conservation in both cases. In addition, we identify several disease-associated variants for
444 which plausible mechanisms had previously been unavailable, but for which allosteric
445 mechanisms provide a plausible rationale.

QWB

446 Unlike the characterization of many other structural features, such as secondary
447 structure assignment, residue burial, protein-protein interaction interfaces, disorder, and
448 even stability, allostery inherently manifests in the context of dynamic behavior. It is only
449 by considering protein motions and changes in these motions can a fuller understanding
450 of allosteric regulation be realized. As such, MD and NMR are some of the most
451 common means of studying allostery and dynamic behavior (Kornev and Taylor, 2015).
452 However, these methods have limitations when studying large and diverse protein
453 datasets. MD is computationally expensive and impractical when studying large numbers
454 of proteins. NMR structure determination is extremely labor-intensive and better suited to
455 certain classes of structures or dynamics. In addition, NMR structures constitute a
456 relatively small fraction of structures currently available.

457 Despite these limitations in MD and NMR, allosteric mechanisms and signaling
458 pathways may be conserved across many different but related proteins, suggesting that
459 such computationally- or labor-intensive approaches for all proteins may not be entirely
460 essential. Flock et al. have carefully demonstrated that the allosteric mechanisms
461 responsible for regulating G proteins through GPCRs tend to be conserved (Flock et al.,
462 2015). If allosteric mechanisms are similarly shared within other protein families, a

463 detailed analysis with methods such as MD or NMR on one member of a family may help
464 to elucidate the allosteric behavior for other members. Nevertheless, the degree to which
465 these mechanisms are indeed conserved within other groups of proteins is currently
466 unclear, so homology-based predictions of allosteric mechanisms are still not readily
467 available.

468 Investigations into representative families have also been enlightening in other
469 contexts. In one of the early studies employing network analysis, del Sol et al. conduct a
470 detailed study of several allosteric protein families (including GPCRs) to demonstrate
471 that residues important for maintaining the integrity of short paths within residue contact
472 networks are essential to enabling signal transmission between distant sites (del Sol et al,
473 2006). Notably, many of the key sites identified correspond to residues that had been
474 experimentally determined to be important for allostery. Another notable result in the
475 same work is that these key residues may become redistributed when the protein
476 undergoes conformational change, thereby changing optimal communication routes in
477 different conformations as a means of conferring different regulatory properties.

478 There are several notable implications of our database-scale analysis. Relative to
479 sequence data, allostery and dynamic behavior are far more difficult to evaluate on a
480 large scale. The framework described here enables one to evaluate dynamic behavior in a
481 systemized and efficient way across many proteins, while simultaneously capturing
482 residues on both the surface and within the interior. That this pipeline can be applied in a
483 high-throughput manner enables the investigation of system-wide phenomena, such as
484 the roles of potential allosteric hotspots in protein-protein interaction networks.
485 Knowledge of such sites across many proteins may also be used to identify the best

DECLAN CLARKE 12/8/15 3:18 PM

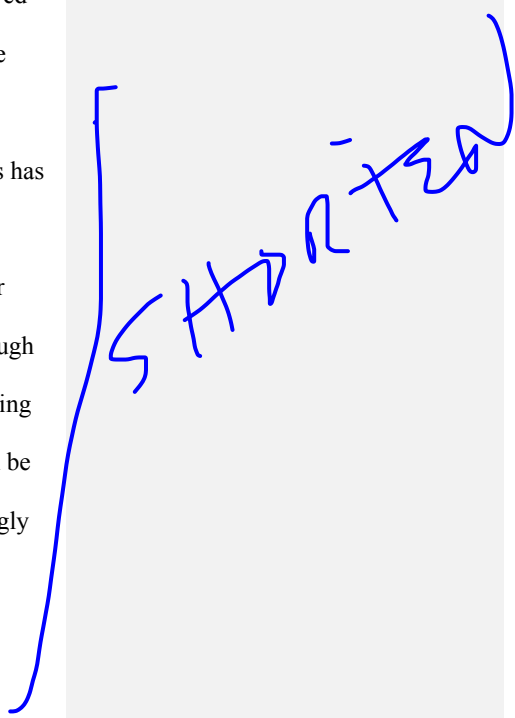
Deleted: Thus, there

487 proteins and protein regions for which drugs should be engineered, as well as instances in
488 which specific sequence variants are likely to have the greatest impact.

489 We emphasize that it is only by applying this framework over a database of many
490 proteins can one search for significant disparities in conservation between sites believed
491 to be important in allostery and the rest of the protein. Such general trends may not be
492 apparent when studying a small number or specific classes of proteins. To our
493 knowledge, this is the first study in which the conservation of potential allosteric sites has
494 been measured across a large database of proteins.

495 The ability to leverage our framework in a high-throughput manner also better
496 enables one to match structural features with the high-throughput data generated through
497 deep sequencing. Full human genomes and exomes are being sequenced at an increasing
498 pace, thereby providing an unprecedented window into conservation patterns that can be
499 human-specific or active over short evolutionary timescales. These patterns increasingly
500 serve as detailed signatures of selective constraints which may not only be missing in
501 cross-species comparisons, but are also sometimes difficult to rationalize using static
502 representations of protein structures alone.

503 We anticipate that, within the next decade, deep sequencing will enable structural
504 biologists to study evolutionary conservation using sequenced human exomes just as
505 routinely as cross-species alignments. Furthermore, intra-species metrics for conservation
506 provide added value in that the confounding factors of cross-species comparisons are
507 removed: different organisms evolve in different cellular and evolutionary contexts, and
508 it can be difficult to decouple these different effects from one another. Cross-species
509 metrics of protein conservation entail comparisons between proteins that may be very



SHARITZ

510 different in structure and function. Sequence-variable regions across species may not be
511 conserved, but nevertheless impart essential functionality. Intra-species comparisons,
512 however, can often provide a more direct and sensitive evaluation of constraint. In
513 addition, intra-species selective constraints are particularly relevant in the context of
514 human disease. Finally, we anticipate that our newly developed software tool will prove
515 to be of great value in enabling investigators to study allostery in diverse contexts.

BETTER LINK

517 METHODS

518 An overview of the framework for finding surface- and interior-critical residues is
519 given in Figure 1. Figure [S3](#) provides a schematic of our pipeline for identifying
520 alternative conformations throughout the PDB. Cross-species conservation scores were
521 analyzed in those PDBs for which full ConSurf files are available through the ConSurf
522 server. 1000 Genomes SNVs were taken from the Phase 3 release, and ExAC SNVs were
523 downloaded in May 2015. Further details on all [protocols](#) are provided in SI Methods.

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: S9

524

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: methods

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: -

525 ACKNOWLEDGMENTS

526 DC acknowledges the support of the NIH Predoctoral Program in Biophysics (T32
527 GM008283-24). We thank Simon Mitternacht for sharing the original source code for
528 binding leverage calculations, as well as Koon-Kiu Yan for helpful discussions and
529 feedback. The authors would like to thank the Exome Aggregation Consortium and the

533 groups that provided exome variant data for comparison. A full list of contributing groups
534 can be found at <http://exac.broadinstitute.org/about>

535

536

537 REFERENCES

- 538 Adzhubei, I. Schmidt, S., Peshkin, L., Ramensky, V.E., Gerasimova, A., Bork, P.,
539 Kondrashov, A.S., and Sunyaev, S.R. (2010). A method and server for predicting
540 damaging missense mutations. *Nat. Methods.* 7, 248–249
- 541 Arnlund, D., Johansson, L.C., Wickstrand, C., Barty, A., Williams, G.J., Malmerberg, E.,
542 Davidsson, J., Milathianaki, D., DePonte, D.P., Shoeman, R.L., et al. (2014).
543 Visualizing a protein quake with time-resolved X-ray scattering at a free-electron
544 laser. *Nat. Methods.* 11, 923–6.
- 545 Atilgan, A.R., Durell, S.R., Jernigan, R.L., Demirel, M.C., Keskin, O., and Bahar, I.
546 (2001). Anisotropy of Fluctuation Dynamics of Proteins with an Elastic Network
547 Model. *Biophys. J.* 80, 505–515.
- 548 Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H.,
549 Shindyalov, I.N. and Bourne, P.E. (2000). The Protein Data Bank. *Nucleic Acids*
550 *Res.* 28, 235–242.
- 551 Calabrese, R., Capriotti, E., Fariselli, P., Martelli, P.L. and Casadio, R. (2009).
552 Functional annotations improve the predictive score of human disease-related
553 mutations in proteins. *Hum. Mutat.* 30, 1237–1244.
- 554 Exome Aggregation Consortium (ExAC). (2015) Cambridge, MA.

555 <http://exac.broadinstitute.org>.

556 Capra, J.A., Laskowski, R.A., Thornton, J.M., Singh, M. and Funkhouser, T.A. (2009).

557 Predicting protein ligand binding sites by combining evolutionary sequence

558 conservation and 3D structure. *PLoS Comput. Biol.* 5, e1000585.

559 Chennubhotla, C. and Bahar, I. (2006). Markov propagation of allosteric effects in

560 biomolecular systems: application to GroEL–GroES. *Mol. Syst. Biol.* 2.

561 [del Sol, A., Fujihashi, H., Amoros, D., and Nussinov, R. \(2006\). Residues crucial for](#)

562 [maintaining short paths in network communication mediate signaling in proteins.](#)

563 [Mol. Syst. Biol.](#) 2(1).

564 Csermely, P., Korecsmáros, T., Kiss, H.J.M., London, G., and Nussinov, R. (2013).

565 Structure and dynamics of molecular networks: A novel paradigm of drug discovery.

566 *Pharmacol. Ther.* 138, 333–408.

567 Flock, T., Ravarani, C.N.J., Sun, D., Venkatakrishnan, A.J., Kayikci, M., Tate, C.G.,

568 Veprintsev, D.B. and Babu, M.M. (2015). Universal allosteric mechanism for Gα

569 activation by GPCRs. *Nature* 524, 173–179.

570 Fuglebakk, E., Tiwari, S.P., and Reuter, N. (2015). Comparing the intrinsic dynamics of

571 multiple protein structures using elastic network models. *Biochim. Biophys. Acta -*

572 *Gen. Subj.* 1850, 911–922.

573 Gasper, P.M., Fuglestad, B., Komives, E.A., Markwick, P.R.L., and McCammon, J.A.

574 (2012). Allosteric networks in thrombin distinguish procoagulant vs. anticoagulant

575 activities. *Proc. Natl. Acad. Sci. U. S. A.* 109, 21216–22.

576 [Ghosh, A., and Vishveshwara, S. \(2008\). Variations in Clique and Community Patterns in](#)

577 [Protein Structures during Allosteric Communication: Investigation of Dynamically](#)

DECLAN CLARKE 12/8/15 3:18 PM

Deleted:

579 | [Equilibrated Structures of Methionyl tRNA Synthetase Complexes. *Biochemistry*.](#)
580 | [47, 11398-11407.](#)

581 | Girvan, M., Girvan, M., Newman, M.E.J., and Newman, M.E.J. (2002). Community
582 | structure in social and biological networks. *Proc. Natl. Acad. Sci. U. S. A.* 99, 7821–
583 | 7826.

584 | Halabi, N., Rivoire, O., Leibler, S., and Ranganathan, R. (2009). Protein Sectors:
585 | Evolutionary Units of Three-Dimensional Structure. *Cell* 138, 774–786.

586 | Huang, B., and Schroeder, M. (2006). LIGSITEcsc: predicting ligand binding sites using
587 | the Connolly surface and degree of conservation. *BMC Struct. Biol.* 6, 19.

588 | Huang, W., Lu, S., Huang, Z., Liu, X., Mou, L., Luo, Y., Zhao, Y., Liu, Y., Chen, Z.,
589 | Hou, T., et al. (2013). AlloSite: A method for predicting allosteric sites.
590 | *Bioinformatics* 29, 2357–2359.

591 | Khurana, E., Fu, Y., Colonna, V., Mu, X.J., Kang, H.M., Lappalainen, T., Sboner, A.,
592 | Lochovsky, L., Chen, J., Harmanci, A., et al. (2013). Integrative Annotation of
593 | Variants from 1092 Humans: Application to Cancer Genomics. *Science*. 342,
594 | 1235587–1235587.

595 | Kornev, A.P. and Taylor, S.S. (2015). Dynamics-Driven Allostery in Protein Kinases.
596 | *Trends Biochem. Sci.* xx, 1–20.

597 | Landrum, M.J., Lee, J.M., Riley, G.R., Jang, W., Rubinstein, W.S., Church, D.M., and
598 | Maglott, D.R. (2014). ClinVar: public archive of relationships among sequence
599 | variation and human phenotype. *Nucleic Acids Res.* 42, D980–5.

600 | Lee, J., Natarajan, M., Nashine, V.C., Socolich, M., Vo, T., Russ, W.P., Benkovic, S.J.,
601 | and Ranganathan, R. (2008). Surface Sites for Engineering Allosteric Control in

602 Proteins. *Science* 322, 438-442.

603 Lockless, S.W., Ranganathan, R., Kukic, P., Mirabello, C., Tradigo, G., Walsh, I., Veltri,
604 P., Pollastri, G., Socolich, M., Lockless, S.W., et al. (1999). Evolutionarily
605 conserved pathways of energetic connectivity in protein families. *BMC*
606 *Bioinformatics* 15, 295–299.

607 McVean, G.A., Altshuler (Co-Chair), D.M., Durbin (Co-Chair), R.M., Abecasis, G.R.,
608 Bentley, D.R., Chakravarti, A., Clark, A.G., Donnelly, P., Eichler, E.E., Flicek, P.,
609 et al. (2012). An integrated map of genetic variation from 1,092 human genomes.
610 *Nature* 491, 56–65.

611 Ming, D. and Wall, M.E. (2005). Quantifying allosteric effects in proteins. *Proteins* 59,
612 697–707.

613 Mitternacht, S. and Berezovsky, I.N. (2011). Binding leverage as a molecular basis for
614 allosteric regulation. *PLoS Comput. Biol.* 7, e1002148.

615 Miyashita, O., Onuchic, J.N., and Wolynes, P.G. (2003). Nonlinear elasticity, protein
616 quakes, and the energy landscapes of functional transitions in proteins. *Proc. Natl.*
617 *Acad. Sci.* 100, 12570–12575.

618 Ng, P.C. and Henikoff, S. (2001). Predicting Deleterious Amino Acid Substitutions.
619 *Genome Res.* 11, 863–874.

620 Panjkovich, A. and Daura, X. (2012). Exploiting protein flexibility to predict the location
621 of allosteric sites. *BMC Bioinformatics* 13, 273.

622 Panjkovich, A. and Daura, X. (2010). Assessing the structural conservation of protein
623 pockets to study functional and allosteric sites: implications for drug discovery.
624 *BMC Struct. Biol.* 10, 9.

625 Reynolds, K.A., McLaughlin, R.N., and Ranganathan, R. (2011). Hot Spots for Allosteric
626 Regulation on Protein Surfaces. *Cell* *147*, 1564–1575.

627 Rivalta, I., Sultan, M.M., Lee, N.-S., Manley, G. a., Loria, J.P., and Batista, V.S. (2012).
628 PNAS Plus: Allosteric pathways in imidazole glycerol phosphate synthase. *Proc.*
629 *Natl. Acad. Sci.* *109*, E1428–E1436.

630 Rodgers, T.L., Townsend, P.D., Burnell, D., Jones, M.L., Richards, S.A., McLeish,
631 T.C.B., Pohl, E., Wilson, M.R., and Cann, M.J. (2013). Modulation of Global Low-
632 Frequency Motions Underlies Allosteric Regulation: Demonstration in CRP/FNR
633 Family Transcription Factors. *PLoS Biol.* *11*, e1001651.

634 Rousseau, F. and Schymkowitz, J. (2005). A systems biology perspective on protein
635 structural dynamics and signal transduction. *Curr. Opin. Struct. Biol.* *15*, 23–30.

636 Sethi, A., Eargle, J., Black, A.A., and Luthey-Schulten, Z. (2009). Dynamical networks
637 in tRNA:protein complexes. *Proc. Natl. Acad. Sci. U. S. A.* *106*, 6620–5.

638 Shulman, A.I., Larson, C., Mangelsdorf, D.J., and Ranganathan, R. (2004). Structural
639 determinants of allosteric ligand activation in RXR heterodimers. *Cell* *116*, 417–
640 429.

641 Stenson, P.D., Mort, M., Ball, E. V., Shaw, K., Phillips, A.D., and Cooper, D.N. (2014).
642 The Human Gene Mutation Database: building a comprehensive mutation repository
643 for clinical and molecular genetics, diagnostic testing and personalized genomic
644 medicine. *Hum. Genet.* *133*, 1–9.

645 Süel, G.M., Lockless, S.W., Wall, M.A., and Ranganathan, R. (2003). Evolutionarily
646 conserved networks of residues mediate allosteric communication in proteins. *Nat.*
647 *Struct. Biol.* *10*, 59–69.

648 Tirion, M.M. (1996). Large Amplitude Elastic Motions in Proteins from a Single-
649 Parameter, Atomic Analysis. *Phys. Rev. Lett.* 77, 1905–1908.
650 Tsai, C., Ma, B. and Nussinov, R. (1999). Folding and binding cascades: Shifts in energy
651 landscapes. *Proc. Natl. Acad. Sci. U. S. A.* 96, 9970–9972.
652 Vanwart, A.T., Eargle, J., Luthey-Schulten, Z., and Amaro, R.E. (2012). Exploring
653 residue component contributions to dynamical network models of allostery. *J.*
654 *Chem. Theory Comput.* 8, 2949–2961.
655 Yang, L.W. and Bahar, I. (2005). Coupling between catalytic site and collective
656 dynamics: A requirement for mechanochemical activity of enzymes. *Structure* 13,
657 893–904.

658
659

660 CAPTIONS

661 **Figure 1. Schematic overviews of methods for finding surface- and interior-critical**
662 **residues.** (A) A simulated ligand probes the protein surface in a series of Monte Carlo
663 simulations (top-left). The cavities identified may be such that occlusion by the ligand
664 strongly interferes with conformational change (top-right; such a site is likely to be
665 identified as surface-critical, in red), or they may have little effect on conformational
666 change, as in the case of shallow pockets (bottom-left) or pockets in which large-scale
667 motions do not drastically affect pocket volume (bottom-right). (B) Interior-critical
668 residues are identified by weighting residue-residue contacts (edges) on the basis of
669 correlated motions, and then identifying communities within the weighted network.

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Hyperlink, No underline, Font color: Auto

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: FIGURE

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: affect

672 Residues involved in the highest-betweenness interactions between communities (in red)
673 are selected as interior-critical residues.

674

675 **Figure 2. Summary statistics for surface-critical sites.** The distributions of the
676 numbers of surface-critical sites per domain and per complex are given in (A) and (B),
677 respectively. Panel (C) gives the distributions of the number of surface-critical sites per
678 complex without thresholding. Complexes are taken from the the PDB biological
679 assembly files. Without applying thresholds to the list of ranked surface-critical sites, the
680 protein is often covered with an excess of identified critical sites.

681

682 **Figure 3: STRESS web server front page, running times, and server architecture.**

683 (A) The server enables users to either provide PDB IDs or to upload their own PDB files
684 for proteins of interest. Users may opt to identify surface-critical residues, interior-critical
685 residues, or both. (B) Running times are shown for systems of various sizes. Shown in
686 red are the running times without optimizing for speed. Performing local searching
687 supported with hashing and implementing additional algorithmic optimizations for
688 computational efficiency reduce running times considerably (in green), relative to a more
689 naïve approach without optimization (in red). (C) The same data is represented as a log-
690 log plot. The slopes of these two approaches demonstrate that our algorithm reduces the
691 computational complexity by an order of magnitude. Our speed-optimized algorithm
692 scales at $O(n^{1.3})$, where n is the number of residues. (D) A thin front-end server handles
693 incoming user requests, and more powerful back-end servers perform the heavier
694 algorithmic calculations. The back-end servers are dynamically scalable, making them

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: on database-wide analyses

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: The distributions of the number of conformations (i.e., "K") for domains and chains are given in (C) and (D), respectively. Only proteins for which K exceeds 1 (for chains) are included in our dataset of multiple conformations. (E) Distinct proteins in our dataset within the context of high-quality X-ray structures in the PDB that we structurally aligned. A set of distinct proteins is such that no pair shares more than 90% sequence identity

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Font:Not Bold

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Tabs: 1.44", Left

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: Figure 3.

708 [capable of handling wide fluctuations in user demand. Amazon's Simple Queue Service](#)
709 [is used to coordinate between user requests at the front end and the back-end compute](#)
710 [nodes: when the front-end server receives a request, it adds the job to the queue, and](#)
711 [back-end servers pull that job from the queue when ready. Source code is available](#)
712 [through Github \(github.com/gersteinlab/STRESS\).](#)

714 **Figure 4. Multiple metrics and datasets reveal that critical residues tend to be**
715 **conserved.** Surface- and interior-critical residues (red) in phosphofructokinase (PDB
716 3PFK) are given in (A) and (E), respectively. Distributions of cross-species conservation
717 scores, 1000 Genomes SNV DAF averages, and ExAC SNV MAF averages for surface-
718 and non-critical residue sets are given in (B), (C), and (D), respectively. The same
719 distributions corresponding to interior- and non-critical residue sets are given in (F), (G),
720 and (H), respectively. In (C), means for surface- and non-critical sets are 9.10×10^{-4} and
721 8.34×10^{-4} , respectively ($p=0.309$); corresponding means in (D) are 4.09×10^{-4} and 2.26×10^{-4} ,
722 respectively ($p=1.49 \times 10^{-3}$). In (G), means for interior- and non-critical sets are 2.82×10^{-4} and
723 3.12×10^{-3} , respectively ($p=1.80 \times 10^{-5}$); corresponding means in (H) are 3.08×10^{-5} and 3.27×10^{-5} ,
724 respectively ($p=7.98 \times 10^{-9}$). Statistics for panels (B) and (F) are given in the main text.
725 $N = 421, 32, 84, 517, 31,$ and 90 structures for panels B, C, D, F, G, and H, respectively.
726 P-values are based on Wilcoxon-rank sum tests. See SI Methods for further details.

727
728 **Figure 5: Critical residues are shown to be more conserved, as measured by the**
729 **fraction of rare alleles. Protein regions with high fractions of rare variants are believed**
730 **to be more sensitive to sequence variants than other regions, thereby explaining why such**

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Tabs: 1.44", Left

DECLAN CLARKE 12/8/15 3:18 PM

Deleted: 4

732 variants occur infrequently in the population. Panels (A) and (C) show distributions for
733 rare (low DAF) non-synonymous SNVs (taken from the 1000 Genomes dataset) in which
734 the critical residues are defined to be the surface-critical (A) and interior-critical (C)
735 residues. Panels (B) and (D) show distributions for rare (low MAF) non-synonymous
736 SNVs (taken from the ExAC dataset) in which the critical residues are defined to be the
737 surface-critical (B) and interior-critical (D) residues. For varying thresholds to define
738 rarity, there are more structures in which the fraction of rare variants is higher in critical
739 residues than in non-critical residues. Cases in which the fraction is equal in both
740 categories are not shown. We consider all structures such that at least one critical and at
741 least one non-critical residue are hit by a non-synonymous SNV. Panels (A), (B), (C), and
742 (D) represent data from 31, 90, 32, and 84 structures, respectively.

743

744 **Figure 6: Modeling protein conformational change through a direct use of crystal**
745 **structures from alternative conformations using absolute conformational transitions**
746 **(ACT).** (A) Distributions (155 structures) of the mean conservation scores on surface-
747 critical (red) and non-critical residues with the same degree of burial (blue). (B)
748 Distributions (159 structures) of the mean conservation scores for interior-critical (red)
749 and non-critical residues with the same degree of burial (blue). Mean values are given in
750 parentheses. Results for single-chain proteins are shown, and p-values were calculated
751 using a Wilcoxon rank sum test.

752

753 **Figure 7. Potential allosteric residues add a layer of annotation to structures in the**
754 **context of disease-associated SNVs.** The structure shown (A) is that of the fibroblast

755 growth-factor receptor (FGFR) in VMD Surf rendering, with HGMD SNVs shown in
756 orange, bound to FGF2, in ribbon rendering (PDB 1IIL). (B) A linear representation of
757 structural annotation for FGFR. Dotted lines highlight loci which correspond to HGMD
758 sites that coincide with critical residues, but for which other annotations fail to coincide.
759 Deeply-buried residues are defined to be those that exhibit a relative solvent-exposed
760 surface area of 5% or less, and binding site residues are defined as those for which at
761 least one heavy atom falls within 4.5 Angstroms of any heavy atom in the binding partner
762 (heparin-binding growth factor 2). The loci of PTM sites were taken from UniProt
763 (accession P21802).

764

765 **Table 1: Statistics on the surfaces of *apo* structures within the canonical set of**
766 **proteins**

767 For each *apo* structure within the canonical set of proteins, statistics relating surface-
768 critical sites to known ligand-binding sites are reported. The surface of a given structure
769 is defined to be the set of all residues that have a relative solvent accessibility of at least
770 50%, where relative solvent accessibility is evaluated using all heavy atoms in both the
771 main-chain and side-chain of a given residue. Mean values are given in the bottom row.
772 NACCESS is used to calculate relative solvent accessibility (Hubbard and Thornton,
773 1993) . *Column 1*: PDB IDs for each structure; *Column 2*: among these surface residues,
774 the fraction that constitute surface-critical residues; *Column 3*: among surface residues,
775 the fraction that constitute known ligand-binding residues (known ligand-binding
776 residues are taken to be those within 4.5 Angstroms of the ligand in the *holo* structure;
777 Table S1); *Column 4*: the Jaccard similarity between the sets of residues represented in

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Font:Not Bold

DECLAN CLARKE 12/8/15 3:18 PM

Formatted: Tabs: 1.44", Left

778 columns 2 and 3 (i.e., surface-critical and known-ligand binding residues), where values
779 given in parentheses represent the expected Jaccard similarity, given a null model in
780 which surface-critical and ligand-binding residues are randomly distributed throughout
781 the surface (for each structure, 10,000 simulations are performed to produce random
782 distributions, and the expected values reported here constitute the mean Jaccard similarity
783 among the 10,000 simulations for each structure); *Column 5*: the number of distinct
784 surface-critical sites identified in each structure; *Column 6*: the number of known ligand-
785 binding sites in each structure; *Column 7*: the number of known ligand-binding sites
786 which are positively identified within the set of surface-critical sites, where a positive
787 match occurs if a majority of the residues in a surface-critical site coincide with the
788 known ligand-binding site; *Column 8*: The fraction of ligand-binding sites captured is
789 simply the ratio of the values in column 7 to those in column 6.
790

Page 6: [1] Deleted

DECLAN CLARKE

12/8/15 3:18 PM

Page 7: [2] Deleted

DECLAN CLARKE

12/8/15 3:18 PM

Page 10: [3] Deleted

DECLAN CLARKE

12/8/15 3:18 PM