1 **Title:**
2 Identifying allosteric hotspots with dynamics: application to inter- and intra-species
3 conservation
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 [a] Department of Chemistry, Yale University, 260/266 Whitney Avenue PO Box 208114,
10 New Haven, CT 06520 USA
11 [b] Program in Computational Biology and Bioinformatics, Yale University, 260/266
12 Whitney Avenue PO Box 208114, New Haven, CT 06520, USA
13 [c] Department of Molecular Biophysics and Biochemistry, Yale University, 260/266
14 Whitney Avenue PO Box 208114, New Haven, CT 06520, USA
15 [d] Department of Computer Science, Yale University, 260/266 Whitney Avenue PO Box
16 208114, New Haven, CT 06520, USA
17 [e] Yale College, 260/266 Whitney Avenue PO Box 208114, New Haven, CT 06520, USA
18 [f] Integrated 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)
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

# ABSTRACT

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

absent. Allostery may often provide the missing explanatory link. We use models of

protein conformational change to identify allosteric residues by finding essential surface

cavities and 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, our analysis is computationally

fast, thereby allowing us to run it across the PDB and to evaluate general properties of

predicted allosteric residues. We find that these tend to be conserved over diverse

evolutionary time scales. Finally, we highlight examples of allosteric residues that help

explain poorly understood disease-associated variants.

# **INTRODUCTION**

69        The ability to sequence large numbers of human genomes is providing a much

70    deeper view into protein evolution than previously possible. When trying to understand

71    the evolutionary pressures on a given protein, structural biologists now have at their

72    disposal an unprecedented breadth of data regarding patterns of conservation, both across

73    species and amongst humans. As such, there are greater opportunities to take an

74    integrated view of the context in which a protein and its residues function. This view

75    necessarily includes structural constraints such as residue packing, protein-protein

76    interactions, and stability. However, deep sequencing is unearthing a class of conserved

77    residues on which no obvious structural constraints appear to be acting. The missing link

78    in understanding these regions may be provided by studying the protein's dynamic

79    behavior through the lens of the distinct functional and conformational states within an

80    ensemble.

81        The underlying energetic landscape responsible for the relative distributions of

82    alternative conformations is dynamic in nature: allosteric signals or other external

83    changes may reconfigure and reshape the landscape, thereby shifting the relative

84    populations of states within an ensemble (Tsai *et al*., 1999). Landscape theory thus

85    provides the conceptual underpinnings necessary to describe how proteins change

86    behavior and shape under changing conditions. A primary driving force behind the

87    evolution of these landscapes is the need to efficiently regulate activity in response to

88    changing cellular contexts, thereby making allostery and conformational change essential

89    components of protein evolution.

90      Given the importance of allosteric regulation, as well as its role in imparting

91      efficient functionality, several methods have been devised for the identification of likely

92      allosteric residues. Conservation itself has been used, either in the context of conserved

93      residues (Panjkovich and Daura, 2012), networks of co-evolving residues (Halabi *et al.*,

94      2009; Lee *et al.*, 2008; Lockless *et al.*, 1999; Reynolds *et al.*, 2011; Shulman *et al.*, 2004;

95      Süel *et al.*, 2003), or local conservation in structure (Panjkovich and Daura, 2010). In

96      related studies, both conservation and geometric-based searches for allosteric sites have

97      been successfully applied to several systems (Capra *et al.*, 2009).

98      The concept of 'protein quakes' has been introduced to explain local

99      conformational changes that are essential for global conformation transitions of

100    functional importance (Ansari *et al.*, 1985; Miyashita *et al.*, 2003). These local changes

101    cause strain within the protein that is relieved by subsequent relaxations (which are also

102    termed functionally important motions) that terminate when the protein reaches the

103    second equilibrium state. Such local perturbations often end with large conformational

104    changes at the focal points of allosteric regulation, and these motions may be identified in

105    a number of ways, including modified normal modes analysis (Miyashita *et al.*, 2003) or

106    time-resolved X-ray scattering (Arnlund *et al.*, 2014).

107    In addition to conservation and geometry, protein dynamics have also been used

108    to predict allosteric residues. Normal modes analysis has been used to examine the extent

109    to which bound ligands interfere with low-frequency motions, thereby identifying

110    potentially important residues at the surface (Ming and Wall, 2005; Mitternacht and

111    Berezovsky, 2011; Panjkovich and Daura, 2012). Normal modes have also been used by

112    the Bahar group to identify important subunits that act in a coherent manner for specific

4

113    proteins (Chennubhotla and Bahar, 2006; Yang and Bahar, 2005). Rodgers *et al.* have

114    applied normal modes to identify key residues in CRP/FNR transcription factors

115    (Rodgers *et al.*, 2013).

116           With the objective of identifying allosteric residues within the interior, molecular

117    dynamics (MD) simulations and network analyses have been used to identify residues

118    that may function as internal allosteric bottlenecks (Csermely *et al.*, 2013; Gasper *et al.*,

119    2012; Rousseau and Schymkowitz, 2005; Sethi *et al.*, 2009; Vanwart *et al.*, 2012). Ghosh

120    *et al.* (2008) have taken a novel approach of combining MD and network principles to

121    characterize allosterically important communication between domains in methionyl

122    tRNA synthetase. In conjunction with NMR, Rivalta *et al.* have use MD and network

123    analysis to identify important regions in imidazole glycerol phosphate synthase (Rivalta

124    *et al.*, 2012).

125           Though having provided valuable insights, many of these approaches have been

126    limited in terms of scale (the numbers of proteins which may feasibly be investigated),

127    computational demands, or the class of residues to which the method is tailored (surface

128    or interior). Here, we use models of protein conformational change to identify both

129    surface and interior residues that may act as essential allosteric hotspots in a

130    computationally tractable manner, thereby enabling high-throughput analysis. This

131    framework directly incorporates information regarding 3D protein structure and

132    dynamics, and it can be applied on a PDB-wide scale to proteins that exhibit

133    conformational change. Throughout the PDB (Berman *et al.*, 2000), the residues

134    identified tend to be conserved both across species and amongst humans, and they may

135    help to elucidate many of the otherwise poorly understood regions in proteins. In a

136 similar vein, several of our identified sites correspond to human disease loci for which no

137 clear mechanism for pathogenesis had previously been proposed. Finally, we make the

138 software associated with this framework (termed STRESS, for STRucturally-identified

139 ESSential residues) publically available through a tool to enable users to submit their

140 own structures for analysis.

141

# 142 RESULTS

## 143 Identifying Potential Allosteric Residues

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

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

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

147 act to transmit such signals. We use models of protein conformational change to identify

148 both classes of residues (Figure 1). Throughout, we term these potential allosteric

149 residues at the surface and interior "surface-critical" and "interior-critical" residues,

150 respectively.

151       In order to gauge the effectiveness of our approach, we identified and analyzed

152 critical residues within a set of 12 well-studied canonical systems (see Figure S1, as well

153 as Table S1 for rationale). We then apply this protocol on a large scale across hundreds of

154 proteins for which crystal structures of alternative conformations are available.

155

156 **Identifying Surface-Critical Residues**

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

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

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

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

161    degree to which cavity occlusion by the ligand disrupts large-scale conformational

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

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

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

165    unaffected (Figure 1A, bottom-right) earn lower scores. Further details are provided in SI

166    Methods section 3.1-a.

167        This approach is a modified version of the binding leverage framework

168    introduced by Mitternacht and Berezovsky (Mitternacht and Berezovsky, 2011). The

169    main modifications implemented here include the use of heavy atoms in the protein

170    during the Monte Carlo search, in addition to an automated means of thresholding the list

171    of ranked scores. These modifications were implemented to provide a more selective set

172    of sites; without them, a very large fraction of the protein surface would be occupied by

173    critical sites (Figure S2C). Within our dataset of proteins exhibiting alternative

174    conformations, we find that this modified approach results in an average of ~2 distinct

175    sites per domain (Figure S2A; see Figure S2B for the distribution for distinct sites within

176    entire complexes).

177        Within the canonical set of 12 proteins, we positively identify an average of 56%

178    of the sites known to be directly involved in ligand or substrate binding (see Table 1,

179    Figure S1, and SI Methods section 3.1-a-iv). Some of the sites identified do not directly

180     overlap with known binding regions, but we often find that these "false positives"

181     nevertheless exhibit some degree of overlap with binding sites (Table S2). In addition,

182     those surface-critical sites that do not match known binding sites may nevertheless

183     correspond to latent allosteric regions: even if no known biological function is assigned

184     to such regions, their occlusion may nevertheless disrupt hitherto unfound large-scale

185     motions.

186

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

188         The binding leverage framework described above is intended to capture hotspot

189     regions at the protein surface, but the Monte Carlo search employed is *a priori* excluded

190     from the protein interior. Allosteric residues often act within the protein interior by

191     functioning as essential information flow 'bottlenecks' within the communication

192     pathways between distant regions.

193         To identify such bottleneck residues, the protein is first modeled as a network,

194     wherein residues represent nodes and edges represent contacts between residues (in much

195     the same way that the protein is modeled as a network in constructing anisotropic

196     network models, see below). In this regard, the problem of identifying interior-critical

197     residues is reduced to a problem of identifying nodes that participate in network

198     bottlenecks (see Figure 1B and SI Methods section 3.1-b for details). Briefly, the network

199     edges are first weighted by the degree of strength in the correlated motions of contacting

200     residues: a strong correlation in the motion between contacting residues implies that

201     knowing how one residue moves better enables one to predict the motion of the other,

202     thereby suggesting a strong information flow between the two residues. The weights are

203    used to assign 'effective distances' between connecting nodes, with strong correlations

204    resulting in shorter effective node-node distances.

205        Using the motion-weighted network, "communities" of nodes are identified using

206    the Girvan-Newman formalism (Girvan *et al*., 2002). This formalism entails calculating

207    the betweenness of each edge, where the betweenness of a given edge is defined as the

208    number of shortest paths between all pairs of residues that pass through that edge (each

209    path length is the sum of that path's effective node-node distances assigned in the

210    weighting scheme above). Each community identified is a group of nodes such that each

211    node within the community is highly inter-connected, but loosely connected to other

212    nodes outside the community. Communities are thus densely inter-connected regions

213    within proteins. As tangible examples, the community partitions and the resultant critical

214    residues for the canonical set are given in Figure 2.

215        Those residues that are involved in the highest-betweenness edges between pairs

216    of interacting communities are identified as the interior-critical residues. These residues

217    are essential for information flow between communities, as their removal would result in

218    substantially longer paths between the residues of one community to those of another.

219

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

221        We have made the implementations for finding surface- and interior-critical

222    residues available through a new software tool, STRESS, which may be accessed at

223    stress.molmovdb.org (Figure 3A). Users may submit a PDB file or a PDB ID

224    corresponding to a structure to be analyzed, and the output provided constitutes the set of

225    identified critical residues.

226       Running times are minimized by using a scalable server architecture that runs on

227    the Amazon cloud (Figure 3D). A light front-end server handles incoming user requests,

228    and more powerful back-end servers, which perform the calculations, are automatically

229    and dynamically scalable, thereby ensuring that they can handle varying levels of demand

230    both efficiently and economically. In addition, the algorithmic implementation of our

231    software is highly efficient, thereby obviating the need for long wait times. Relative to a

232    naïve global Monte Carlo search implementation, local searches supported with hashing

233    and additional algorithmic optimizations for computational efficiency reduce running

234    times considerably (Figures 3B and 3C). A typical protein of ~500 residues takes only

235    about 30 minutes on a 2.6GHz CPU.

236

## 237 High-Throughput Identification of Alternative

## 238 Conformations

239       We use a generalized approach to systematically identify instances of alternative

240    conformations throughout the PDB. We first perform multiple structure alignments

241    (MSAs) across sequence-identical structures that are pre-filtered to ensure structural

242    quality. We then use the resultant pairwise RMSD values to infer distinct conformational

243    states (Figure S3; see also SI Methods section 3.2).

244       The distributions of the resultant numbers of conformations for domains and

245    chains are given in Figures S3D and S3E, respectively, and an overview is given in

246    Figure S3F. We note that the alternative conformations identified arise in an extremely

247    diverse set of biological contexts, including conformational transitions that accompany

248  ligand binding, protein-protein or protein-nucleic acid interactions, post-translational

249  modifications, changes in oxidation or oligomerization states, etc. The dataset of

250  alternative conformations identified is provided as a resource in File S1 (see also Figure

251  S3G).

252

## Evaluating Conservation of Critical Residues Using Various Metrics and Sources of Data

255  The large dataset of dynamic proteins culled throughout the PDB, coupled with

256  the high algorithmic efficiency of our critical residue search implementation, provide a

257  means of identifying and evaluating general properties of a large pool of critical residues.

258  In particular, we use a variety of conservation metrics and data sources to measure the

259  inter- and intra-species conservation of the residues within this pool. As discussed below,

260  we find that both surface- (Figures 4A-D) and interior-critical residues (Figures 4E-H)

261  are consistently more conserved than non-critical residues. We emphasize that the

262  signatures of conservation identified not only provide a means of rationalizing many of

263  the otherwise poorly understood regions of proteins, but they also reinforce the functional

264  importance of the residues predicted to be allosteric.

265

266  **Conservation Across Species**

267  When evaluating conservation across species, we find that both surface- and

268  interior-critical residues tend to be significantly more conserved than non-critical residues

269  with the same degree of burial (Figures 4B and 4F, respectively; note that negative

270  conservation scores designate stronger conservation – see SI Methods section 3.3-a).

271

**Leveraging Next-Generation Sequencing to Measure Conservation Amongst**

**Humans**

In addition to measuring inter-species conservation, we have also used fully

sequenced human genomes and exomes to investigate conservation among human

populations, as many constraints may be species-specific and active in more recent

evolutionary history. Commonly used metrics for quantifying intra-species conservation

include minor allele frequency (MAF) and derived allele frequency (DAF). Low MAF or

DAF values are interpreted as signatures of deleteriousness, as purifying selection is

prone to reduce the frequencies of harmful variants (see SI Methods section 3.3-b).

Non-synonymous single-nucleotide variants (SNVs) from the 1000 Genomes

dataset (McVean *et al*., 2012) that intersect surface-critical residues tend to occur at

lower DAF values than do SNVs that intersect non-critical residues (Figure 4C). Though

this difference is not observed to be significant, the significance improves when

examining the shift in DAF distributions, as evaluated with a KS test (p= 0.159, Figure

S4A), and we point out only a limited number of proteins (thirty-two) for which these

1000 Genomes SNVs intersect with surface-critical sites. Furthermore, the long tail

extending to lower DAF values for surface-critical residues may suggest that only a

subset of the residues in our prioritized binding sites is essential. In contrast to surface-

critical residues, however, interior-critical residues intersect 1000 Genomes SNVs with

significantly lower DAF values than do non-critical residues (Figure 4G; see also Figure

S4B).

[[DC2MG(dec18): The paragraph added below was introduced after you

suggested that we discuss the stats issue in the most recent annotated PDF. Another

295

296   When analyzing human polymorphism data, a variety of

297   statistical measures relating SNVs to selective constraint may be calculated, and the

298   results obtained (along with their associated significance levels) are highly dependent on

299   sample size. 1000 Genomes datasets are attractive partially because of their status as a

300   well-established "blue chip" set of variants in human populations. However, given the

301   relatively limited number of proteins that intersect with 1000 Genomes SNVs, we also

302   analyzed the larger dataset provided by the Exome Aggregation Consortium (ExAC)

303   (Exome Aggregation Consortium, 2015). Though this dataset has been released much

304   more recently (and is consequently not yet as well established as 1000 Genomes), ExAC

305   provides sequence data from more than 60,000 individuals, and samples are sequenced at

306   much higher coverage, thereby ensuring better data quality. This larger dataset enables us

307   to more easily examine trends in the data as they relate to critical and non-critical

308   residues.

309         Using MAF as a conservation metric, we performed a similar analysis using this

310   data. MAF distributions for surface- and non-critical residues in the same set of proteins

311   are given in Figure 4D. Although the mean value of the MAF distribution for surface-

312   critical residues is slightly higher than that of non-critical residues, the median for

313   surface-critical residues is substantially lower than that for non-critical residues,

314   demonstrating that the majority of proteins are such that MAF values are lower in

315   surface- than in non-critical residues. In addition, the overall shifts of these distributions

316   also point to a trend of lower MAF values in surface-critical residues (Figure S5A, KS

317   test p=9.49e-2).

318        Interior-critical residues exhibit significantly lower MAF values than do non-

319    critical residues in the same set of proteins. MAF distributions for interior- and non-

320    critical residues are given in Figure 4H (see also Figure S5B).

321        In addition to analyzing overall allele frequency distributions, we also evaluate

322    the *fraction* of rare alleles as a metric for measuring selective pressure. This fraction is

323    defined as the ratio of the number of rare (i.e., low-DAF or low-MAF) non-synonymous

324    SNVs to the number of all non-synonymous SNVs in a given protein annotation (such as

325    all surface-critical residues of the protein, for example; see SI Methods section 3.3-b). A

326    higher fraction is interpreted as a proxy for greater conservation (Khurana *et al*., 2013;

327    Sethi *et al*., 2015). Using variable DAF (MAF) cutoffs to define rarity for 1000 Genomes

328    (ExAC) SNVs, both surface- and interior-critical residues are shown to harbor a higher

329    fraction of rare alleles than do non-critical residues, further suggesting a greater degree of

330    evolutionary constraint on critical residues (See Figure 5).

331

332    **Comparisons Between Different Models of Protein Motions**

333        The identification of surface- and interior-critical residues entails using sets of

334    vectors (on each protein residue) to describe conformational change. Notably, our

335    framework enables one to determine these vectors in multiple ways. Conformational

336    changes may be modeled using vectors connecting residues in crystal structures from

337    alternative conformations. We term this approach "ACT", for "absolute conformational

338    transitions" (see SI Methods section 3.2-c). The crystal structures of such paired

339    conformations may be obtained using the framework discussed above. The protein

340    motions may also be inferred from anisotropic network models (ANMs) (Atilgan *et al*.,

341    2001). ANMs entail modeling interacting residues as nodes linked by flexible springs, in

14

342　a manner similar to elastic network models (Fuglebakk *et al.*, 2015; Tirion, 1996) or

343　normal modes analysis (Figure 1B). ANMs are not only simple and straightforward to

344　apply on a database scale, but unlike using alternative crystal structures, the motion

345　vectors inferred may be generated using a single structure.

346　　　　We find that modeling conformational change using vectors from either ACTs or

347　ANMs gives the same general trends in terms of the disparities in conservation between

348　critical and non-critical residues. Our framework is thus general with respect to how the

349　motion vectors are obtained (see Figure 6 and SI Methods section 3.2-c for further

350　details).

351

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

353　　　　Directly related to conservation is confidence with which an SNV is believed to

354　be disease-associated. SIFT (Ng and Henikoff, 2001) and PolyPhen (Adzhubei *et al.*,

355　2010) are two tools for predicting SNV deleteriousness. ExAC SNVs that intersect

356　critical residues exhibit significantly higher PolyPhen scores relative to non-critical

357　residues, suggesting the potentially higher disease susceptibility at critical residues

358　(Figure S6). Significant disparities were not observed in SIFT scores (Figure S7).

359　　　　Using HGMD (Stenson *et al.*, 2014) and ClinVar (Landrum *et al.*, 2014), we

360　identify proteins with critical residues that coincide with disease-associated SNVs (Figure

361　7A and File S2). Several critical residues coincide with known disease loci for which the

362　mechanism of pathogenicity is otherwise unclear (File S3). The fibroblast growth factor

363　receptor (FGFR) is a case-in-point (Figure 7). SNVs in FGFR have been linked to

364　craniofacial defects. Dotted lines in Figure 7B highlight poorly understood disease SNVs

365　that coincide with critical residues. In addition, we identify Y328 as a surface-critical

366     residue, which coincides with a disease-associated SNV from HGMD, despite the lack of

367     confident predictions of deleteriousness by several widely used tools for predicting

368     disease-associated SNVs, including PolyPhen (Adzhubei *et al*., 2010), SIFT (Ng and

369     Henikoff, 2001), and SNPs&GO (Calabrese *et al*., 2009). Together, these results suggest

370     that the incorporation of surface- and interior-critical residues introduces a valuable layer

371     of annotation to the protein sequence, and may help to explain otherwise poorly

372     understood disease-associated SNVs.

373

374 # DISCUSSION & CONCLUSIONS

375        The same principles of energy landscape theory that dictate protein folding are

376     integral to how proteins explore different conformations once they adopt their fully

377     folded states. These landscapes are shaped not only by the protein sequence itself, but

378     also by extrinsic conditions. Such external factors often regulate protein activity by

379     introducing allosteric-induced changes, which ultimately reflect changes in the shapes

380     and population distributions of the energetic landscape. In this regard, allostery provides

381     an ideal platform from which to study protein behavior in the context of their energetic

382     landscapes. To investigate allosteric regulation, and to simultaneously add an extra layer

383     of annotation to conservation patterns, an integrated framework to identify potential

384     allosteric residues is essential. We introduce a framework to select such residues, using

385     knowledge of conformational change.

386        When applied to many proteins with distinct conformational changes in the PDB,

387     we investigate the conservation of potential allosteric residues in both inter-species and

388     intra-human genomes contexts, and find that these residues tend to exhibit greater

389    conservation in both cases. In addition, we identify several disease-associated variants for

390    which plausible mechanisms had been unknown, but for which allosteric mechanisms

391    provide a reasonable rationale.

392          Unlike the characterization of many other structural features, such as secondary

393    structure assignment, residue burial, protein-protein interaction interfaces, disorder, and

394    even stability, allostery inherently manifests through dynamic behavior. It is only by

395    considering protein motions and changes in these motions can a fuller understanding of

396    allosteric regulation be realized. As such, MD and NMR are some of the most common

397    means of studying allostery and dynamic behavior (Kornev and Taylor, 2015). However,

398    these methods have limitations when studying large and diverse protein datasets. MD is

399    computationally expensive and impractical when studying large numbers of proteins.

400    NMR structure determination is extremely labor-intensive and better suited to certain

401    classes of structures or dynamics. In addition, NMR structures constitute a relatively

402    small fraction of structures currently available.

403          Despite these limitations in MD and NMR, allosteric mechanisms and signaling

404    pathways may be conserved across many different but related proteins within the same

405    family, suggesting that such computationally- or labor-intensive approaches for all

406    proteins may not be entirely essential. Flock *et al*. have carefully demonstrated that the

407    allosteric mechanisms responsible for regulating G proteins through GPCRs tend to be

408    conserved (Flock *et al*., 2015). Investigations into representative families have also been

409    enlightening in other contexts. In one of the early studies employing network analysis,

410    del Sol *et al*. conduct a detailed study of several allosteric protein families (including

411    GPCRs) to demonstrate that residues important for maintaining the integrity of short

412 paths within residue contact networks are essential to enabling signal transmission

413 between distant sites (del Sol *et al.*, 2006). Another notable result in the same work is that

414 these key residues (which match experimental results) may become redistributed when

415 the protein undergoes conformational change, thereby changing optimal communication

416 routes as a means of conferring different regulatory properties.

417 There are several notable implications of our dynamics-based analysis across a

418 database of proteins. Relative to sequence data, allostery and dynamic behavior are far

419 more difficult to evaluate on a large scale. The framework described here enables one to

420 evaluate dynamic behavior in a systemized and efficient way across many proteins, while

421 simultaneously capturing residues on both the surface and within the interior. That this

422 pipeline can be applied in a high-throughput manner enables the investigation of system-

423 wide phenomena, such as the roles of potential allosteric hotspots in protein-protein

424 interaction networks.

425 It is only by analyzing a large dataset of proteins can one investigate general

426 trends in predicted allosteric residues. In addition, the implementation detailed here

427 enables one to match structural features with the high-throughput data generated through

428 deep sequencing initiatives, which are providing an unprecedented window into

429 conservation patterns, many of which may be human-specific.

430 We anticipate that, within the next decade, deep sequencing will enable structural

431 biologists to study evolutionary conservation using sequenced human exomes just as

432 routinely as cross-species alignments. Furthermore, intra-species metrics for conservation

433 provide added value in that the confounding factors of cross-species comparisons are

434 removed: different species evolve in various evolutionary contexts and at different rates,

435   and it can be difficult to decouple these different effects from one another. Cross-species

436   metrics of protein conservation entail comparisons between proteins that may be very

437   different in structure and function. Sequence-variable regions across species may not be

438   conserved, but nevertheless impart essential functionality. Intra-species comparisons,

439   however, can often provide a more direct and sensitive evaluation of constraint.

440           In particular, selective constraints within human populations are particularly

441   relevant to understanding human disease. Formalisms for analyzing large structural and

442   sequence datasets will become increasingly important in the context of human health. We

443   anticipate that the framework and formalisms detailed here, along with the accompanying

444   web tool we have introduced, will help to further motivate future studies along these

445   directions.

446

447   # METHODS

448           An overview of the framework for finding surface- and interior-critical residues is

449   given in Figure 1. Figure S3 provides a schematic of our pipeline for identifying

450   alternative conformations throughout the PDB. Cross-species conservation scores were

451   analyzed in those PDBs for which full ConSurf files are available through the ConSurf

452   server. 1000 Genomes SNVs were taken from the Phase 3 release, and ExAC SNVs were

453   downloaded in May 2015. Further details on all protocols are provided in SI Methods.

454

455   # ACKNOWLEDGMENTS

# REFERENCES

464 Adzhubei, I. Schmidt, S., Peshkin, L., Ramensky, V.E., Gerasimova, A., Bork, P.,

465     Kondrashov, A.S., and Sunyaev, S.R. (2010). A method and server for predicting

466     damaging missense mutations. Nat. Methods. *7*, 248–249

467 Ansari, A., Berendzen, J., Bowne, S., Frauenfelder, H., Iben, I.E.T., Sauke, T.B.,

468     Shyamsunder, E., and Young, R.D. (1985). Protein states and protein quakes. Proc.

469     Natl. Acad. Sci. U.S.A. 82, 5000–5004.

470 Arnlund, D., Johansson, L.C., Wickstrand, C., Barty, A., Williams, G.J., Malmerberg, E.,

471     Davidsson, J., Milathianaki, D., DePonte, D.P., Shoeman, R.L., *et al*. (2014).

472     Visualizing a protein quake with time-resolved X-ray scattering at a free-electron

473     laser. Nat. Methods. *11*, 923–6.

474 Atilgan, A.R., Durell, S.R., Jernigan, R.L., Demirel, M.C., Keskin, O., and Bahar, I.

475     (2001). Anisotropy of Fluctuation Dynamics of Proteins with an Elastic Network

476     Model. Biophys. J. *80*, 505–515.

477 Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H.,

478     Shindyalov, I.N. and Bourne, P.E. (2000). The Protein Data Bank. Nucleic Acids

479     Res. *28*, 235–242.

480     Calabrese, R., Capriotti, E., Fariselli, P., Martelli, P.L. and Casadio, R. (2009).

481         Functional annotations improve the predictive score of human disease-related

482         mutations in proteins. Hum. Mutat. *30*, 1237–1244.

483     Exome Aggregation Consortium, Lek, M., Karczewski, K., Minikel, E., Samocha, K.,

484         Banks, E., Fennell, T., O'Donnell-Luria, A., Ware, J., Hill, A., et al. (2015).

485         Analysis of protein-coding genetic variation in 60,706 humans. bioRxiv. 030338

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

487         Predicting protein ligand binding sites by combining evolutionary sequence

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

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

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

491     del Sol, A., Fujihashi, H., Amoros, D., and Nussinov, R. (2006). Residues crucial for

492         maintaining short paths in network communication mediate signaling in proteins.

493         Mol. Syst. Biol. 2(1).

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

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

496         Pharmacol. Ther. *138*, 333–408.

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

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

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

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

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

502      Gen. Subj. *1850*, 911–922.

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

504      (2012). Allosteric networks in thrombin distinguish procoagulant vs. anticoagulant

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

506    Ghosh, A., and Vishveshwara, S. (2008). Variations in Clique and Community Patterns in

507      Protein Structures during Allosteric Communication: Investigation of Dynamically

508      Equilibrated Structures of Methionyl tRNA Synthetase Complexes. Biochemistry.

509      *47*, 11398-11407.

510    Girvan, M., Girvan, M., Newman, M.E.J., and Newman, M.E.J. (2002). Community

511      structure in social and biological networks. Proc. Natl. Acad. Sci. U. S. A. *99*, 7821–

512      7826.

513    Halabi, N., Rivoire, O., Leibler, S., and Ranganathan, R. (2009). Protein Sectors:

514      Evolutionary Units of Three-Dimensional Structure. Cell *138*, 774–786.

515    Khurana, E., Fu, Y., Colonna, V., Mu, X.J., Kang, H.M., Lappalainen, T., Sboner, A.,

516      Lochovsky, L., Chen, J., Harmanci, A., *et al*. (2013). Integrative Annotation of

517      Variants from 1092 Humans: Application to Cancer Genomics. Science. *342*,

518      1235587–1235587.

519    Kornev, A.P. and Taylor, S.S. (2015). Dynamics-Driven Allostery in Protein Kinases.

520      Trends Biochem. Sci. *xx*, 1–20.

521    Landrum, M.J., Lee, J.M., Riley, G.R., Jang, W., Rubinstein, W.S., Church, D.M., and

522      Maglott, D.R. (2014). ClinVar: public archive of relationships among sequence

523      variation and human phenotype. Nucleic Acids Res. *42*, D980–5.

524    Lee, J., Natarajan, M., Nashine, V.C., Socolich, M., Vo, T., Russ, W.P., Benkovic, S.J.,

525     and Ranganathan, R. (2008). Surface Sites for Engineering Allosteric Control in

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

527     Lockless, S.W., Ranganathan, R., Kukic, P., Mirabello, C., Tradigo, G., Walsh, I., Veltri,

528     P., Pollastri, G., Socolich, M., Lockless, S.W., *et al*. (1999). Evolutionarily

529     conserved pathways of energetic connectivity in protein families. BMC

530     Bioinformatics *15*, 295–299.

531     McVean, G.A., Altshuler (Co-Chair), D.M., Durbin (Co-Chair), R.M., Abecasis, G.R.,

532     Bentley, D.R., Chakravarti, A., Clark, A.G., Donnelly, P., Eichler, E.E., Flicek, P.,

533     *et al*. (2012). An integrated map of genetic variation from 1,092 human genomes.

534     Nature *491*, 56–65.

535     Ming, D. and Wall, M.E. (2005). Quantifying allosteric effects in proteins. Proteins *59*,

536     697–707.

537     Mitternacht, S. and Berezovsky, I.N. (2011). Binding leverage as a molecular basis for

538     allosteric regulation. PLoS Comput. Biol. *7*, e1002148.

539     Miyashita, O., Onuchic, J.N., and Wolynes, P.G. (2003). Nonlinear elasticity, protein

540     quakes, and the energy landscapes of functional transitions in proteins. Proc. Natl.

541     Acad. Sci. *100*, 12570–12575.

542     Ng, P.C. and Henikoff, S. (2001). Predicting Deleterious Amino Acid Substitutions.

543     Genome Res. *11*, 863–874.

544     Panjkovich, A. and Daura, X. (2012). Exploiting protein flexibility to predict the location

545     of allosteric sites. BMC Bioinformatics *13*, 273.

546     Panjkovich, A. and Daura, X. (2010). Assessing the structural conservation of protein

547     pockets to study functional and allosteric sites: implications for drug discovery.

548       BMC Struct. Biol. *10*, 9.

549    Reynolds, K.A., McLaughlin, R.N., and Ranganathan, R. (2011). Hot Spots for Allosteric

550       Regulation on Protein Surfaces. Cell *147*, 1564–1575.

551    Rivalta, I., Sultan, M.M., Lee, N.-S., Manley, G. a., Loria, J.P., and Batista, V.S. (2012).

552       PNAS Plus: Allosteric pathways in imidazole glycerol phosphate synthase. Proc.

553       Natl. Acad. Sci. *109*, E1428–E1436.

554    Rodgers, T.L., Townsend, P.D., Burnell, D., Jones, M.L., Richards, S.A., McLeish,

555       T.C.B., Pohl, E., Wilson, M.R., and Cann, M.J. (2013). Modulation of Global Low-

556       Frequency Motions Underlies Allosteric Regulation: Demonstration in CRP/FNR

557       Family Transcription Factors. PLoS Biol. *11*, e1001651.

558    Rousseau, F. and Schymkowitz, J. (2005). A systems biology perspective on protein

559       structural dynamics and signal transduction. Curr. Opin. Struct. Biol. *15*, 23–30.

560    Sethi, A., Eargle, J., Black, A.A., and Luthey-Schulten, Z. (2009). Dynamical networks

561       in tRNA:protein complexes. Proc. Natl. Acad. Sci. U. S. A. *106*, 6620–5.

562    Sethi, A., Clarke, D., Chen, J., Kumar, S., Galeev, T.R., Regan, L., and Gerstein, M.

563       (2015). Reads meet rotamers: structural biology in the age of deep sequencing. Curr.

564       Opin. Struct. Biol. *35*, 125-134.

565    Shulman, A.I., Larson, C., Mangelsdorf, D.J., and Ranganathan, R. (2004). Structural

566       determinants of allosteric ligand activation in RXR heterodimers. Cell *116*, 417–

567       429.

568    Stenson, P.D., Mort, M., Ball, E. V., Shaw, K., Phillips, A.D., and Cooper, D.N. (2014).

569       The Human Gene Mutation Database: building a comprehensive mutation repository

570       for clinical and molecular genetics, diagnostic testing and personalized genomic

571    medicine. Hum. Genet. *133*, 1–9.

572  Süel, G.M., Lockless, S.W., Wall, M.A., and Ranganathan, R. (2003). Evolutionarily

573    conserved networks of residues mediate allosteric communication in proteins. Nat.

574    Struct. Biol. *10*, 59–69.

575  Tirion, M.M. (1996). Large Amplitude Elastic Motions in Proteins from a Single-

576    Parameter, Atomic Analysis. Phys. Rev. Lett. *77*, 1905–1908.

577  Tsai, C., Ma, B. and Nussinov, R. (1999). Folding and binding cascades: Shifts in energy

578    landscapes. Proc. Natl. Acad. Sci. U. S. A. *96*, 9970–9972.

579  Vanwart, A.T., Eargle, J., Luthey-Schulten, Z., and Amaro, R.E. (2012). Exploring

580    residue component contributions to dynamical network models of allostery. J.

581    Chem. Theory Comput. *8*, 2949–2961.

582  Yang, L.W. and Bahar, I. (2005). Coupling between catalytic site and collective

583    dynamics: A requirement for mechanochemical activity of enzymes. Structure *13*,

584    893–904.

585

586

# CAPTIONS

587

588  **Figure 1.  Schematic overviews of methods for finding surface- and interior-critical**

589  **residues.** (*A*) A simulated ligand probes the protein surface in a series of Monte Carlo

590  simulations (top-left). The cavities identified may be such that occlusion by the ligand

591  strongly interferes with conformational change (top-right; such a site is likely to be

592  identified as surface-critical, in red), or they may have little effect on conformational

593  change, as in the case of shallow pockets (bottom-left) or pockets in which large-scale

594    motions do not drastically affect pocket volume (bottom-right). (*B*) Interior-critical

595    residues are identified by weighting residue-residue contacts (edges) on the basis of

596    correlated motions, and then identifying communities within the weighted network.

597    Residues involved in the highest-betweenness interactions between communities (in red)

598    are selected as interior-critical residues.

599

600    **Figure 2.  Community partitioning for canonical systems.** Different network

601    communities are colored differently, and communities were identified using the

602    dynamical network-based analysis with the GN formalism discussed in the main text and

603    in SI Methods section 3.1-b. Residues shown as spheres are interior-critical residues, and

604    they are colored based on community membership, and black lines connecting pairs of

605    critical residues represent the highest-betweenness edges between the corresponding

606    communities.

607

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

609    (A) The server enables users to either provide PDB IDs or to upload their own PDB files

610    for proteins of interest. Users may opt to identify surface-critical residues, interior-critical

611    residues, or both. (B) Running times are shown for systems of various sizes. Shown in

612    red are the running times without optimizing for speed, and green shows running times

613    with algorithmic optimization. (C) The same data is represented as a log-log plot. The

614    slopes of these two approaches demonstrate that our algorithm reduces the computational

615    complexity by an order of magnitude. Our speed-optimized algorithm scales at $O(n^{1.3})$,

616    where n is the number of residues. (D) A thin front-end server handles incoming user

617  requests, and more powerful back-end servers perform the heavier algorithmic

618  calculations. The back-end servers are dynamically scalable, making them capable of

619  handling wide fluctuations in user demand. Amazon's Simple Queue Service is used to

620  coordinate between user requests at the front end and the back-end compute nodes: when

621  the front-end server receives a request, it adds the job to the queue, and back-end servers

622  pull that job from the queue when ready. Source code is available through Github

623  (github.com/gersteinlab/STRESS).

624

625  **Figure 4. Multiple metrics and datasets reveal that critical residues tend to be**

626  **conserved.** Surface- and interior-critical residues (red) in phosphofructokinase (PDB

627  3PFK) are given in (*A*) and (*E*), respectively. Distributions of cross-species conservation

628  scores, 1000 Genomes SNV DAF averages, and ExAC SNV MAF averages for surface-

629  and non-critical residue sets are given in (*B*), (*C*), and (*D*), respectively. The same

630  distributions corresponding to interior- and non-critical residue sets are given in (F), (G),

631  and (H), respectively. In (B), mean inter-species conservation scores for surface-critical

632  sets are -0.131, whereas non-critical residue sets with the same degree of burial have a

633  mean score of +0.059 (p < 2.2e-16). In (F), mean inter-species conservation scores for

634  interior-critical sets are -0.179, whereas non-critical residue sets with the same degree of

635  burial have a mean score of -0.102 (p=3.67e-11). In (C), means for surface- and non-

636  critical sets are 9.10e-4 and 8.34e-4, respectively (p=0.309); corresponding means in (*D*)

637  are 4.09e-04 and 2.26e-04, respectively (p=1.49e-3). In (*G*), means for interior- and non-

638  critical sets are 2.82e-4 and 3.12e-3, respectively (p=1.80e-05); corresponding means in

639  (*H*) are 3.08e-05 and 3.27e-04, respectively (p=7.98e-09). N = 421, 32, 84, 517, 31, and

640    90 structures for panels B, C, D, F, G, and H, respectively. P-values are based on

641    Wilcoxon-rank sum tests. See SI Methods for further details.

642

643    **Figure 5. Critical residues are shown to be more conserved, as measured by the**

644    **fraction of rare alleles.** Protein regions with high fractions of *rare* variants are believed

645    to be more sensitive to sequence variants than other regions, thereby explaining why such

646    variants occur infrequently in the population. Panels *(A)* and *(C)* show distributions for

647    rare (low DAF) non-synonymous SNVs (taken from the 1000 Genomes dataset) in which

648    the critical residues are defined to be the surface-critical *(A)* and interior-critical *(C)*

649    residues. Panels *(B)* and *(D)* show distributions for rare (low MAF) non-synonymous

650    SNVs (taken from the ExAC dataset) in which the critical residues are defined to be the

651    surface-critical *(B)* and interior-critical *(D)* residues. For varying thresholds to define

652    rarity, there are more structures in which the fraction of rare variants is higher in critical

653    residues than in non-critical residues. Cases in which the fraction is equal in both

654    categories are not shown. We consider all structures such that at least one critical and at

655    least one non-critical residue intersect a non-synonymous SNV. Panels *(A), (B), (C),* and

656    *(D)* represent data from 31, 90, 32, and 84 structures, respectively.

657

658    **Figure 6. Modeling protein conformational change through a direct use of crystal**

659    **structures from alternative conformations using absolute conformational transitions**

660    **(ACT).** *(A)* Distributions (155 structures) of the mean conservation scores on surface-

661    critical (red) and non-critical residues with the same degree of burial (blue). *(B)*

662    Distributions (159 structures) of the mean conservation scores for interior-critical (red)

663 and non-critical residues with the same degree of burial (blue). Mean values are given in

664 parentheses. Results for single-chain proteins are shown, and p-values were calculated

665 using a Wilcoxon rank sum test.

666

667 **Figure 7. Potential allosteric residues add a layer of annotation to structures in the**

668 **context of disease-associated SNVs.** The structure shown (*A*) is that of the fibroblast

669 growth-factor receptor (FGFR) in VMD Surf rendering, with HGMD SNVs shown in

670 orange, bound to FGF2, in ribbon rendering (PDB 1IIL). (*B*) A linear representation of

671 structural annotation for FGFR. Dotted lines highlight loci which correspond to HGMD

672 sites that coincide with critical residues, but for which other annotations fail to coincide.

673 Deeply-buried residues are defined to be those that exhibit a relative solvent-exposed

674 surface area of 5% or less, and binding site residues are defined as those for which at

675 least one heavy atom falls within 4.5 Angstroms of any heavy atom in the binding partner

676 (heparin-binding growth factor 2). The loci of PTM sites were taken from UniProt

677 (accession P21802).

678

679 **Table 1. Statistics on the surfaces of *apo* structures within the canonical set of**

680 **proteins**

681 For each *apo* structure within the canonical set of proteins, statistics relating surface-

682 critical sites to known ligand-binding sites are reported. The surface of a given structure

683 is defined to be the set of all residues that have a relative solvent accessibility of at least

684 50%, where relative solvent accessibility is evaluated using all heavy atoms in both the

685 main-chain and side-chain of a given residue. Mean values are given in the bottom row.

686    NACCESS is used to calculate relative solvent accessibility (Hubbard and Thornton,

687    1993) . *Column 1*: PDB IDs for each structure; *Column 2*: among these surface residues,

688    the fraction that constitute surface-critical residues; *Column 3*: among surface residues,

689    the fraction that constitute known ligand-binding residues (known ligand-binding

690    residues are taken to be those within 4.5 Angstroms of the ligand in the *holo* structure;

691    Table S1); *Column 4*: the Jaccard similarity between the sets of residues represented in

692    columns 2 and 3 (i.e., surface-critical and known-ligand binding residues), where values

693    given in parentheses represent the expected Jaccard similarity, given a null model in

694    which surface-critical and ligand-binding residues are randomly distributed throughout

695    the surface (for each structure, 10,000 simulations are performed to produce random

696    distributions, and the expected values reported here constitute the mean Jaccard similarity

697    among the 10,000 simulations for each structure); *Column 5*: the number of distinct

698    surface-critical sites identified in each structure; *Column 6*: the number of known ligand-

699    binding sites in each structure; *Column 7*: the number of known ligand-binding sites

700    which are positively identified within the set of surface-critical sites, where a positive

701    match occurs if a majority of the residues in a surface-critical site coincide with the

702    known ligand-binding site; *Column 8*: The fraction of ligand-binding sites captured is

703    simply the ratio of the values in column 7 to those in column 6.

704