

 Open access • Posted Content • DOI:10.1101/508788

Leveraging protein dynamics to identify cancer mutational hotspots in 3D-structures

— [Source link](#) 

Sushant Kumar, Declan Clarke, Mark Gerstein

Institutions: Yale University

Published on: 31 Dec 2018 - bioRxiv (Cold Spring Harbor Laboratory)

Related papers:

- [A full-proteome, interaction-specific characterization of mutational hotspots across human cancers](#)
- [Mutational interactions define novel cancer subgroups.](#)
- [CloneSig: Joint inference of intra-tumor heterogeneity and mutational signatures' activity in tumor bulk sequencing data](#)
- [MutSignatures: An R Package for Extraction and Analysis of Cancer Mutational Signatures](#)
- [Network-Based Coverage of Mutational Profiles Reveals Cancer Genes](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/leveraging-protein-dynamics-to-identify-cancer-mutational-2qwf8830wc>

1 **Leveraging protein dynamics to identify cancer mutational hotspots in**
2 **3D-structures**

3
4

5 Sushant Kumar^{1,2}, Declan Clarke^{1,2}, Mark B. Gerstein^{1,2,3*}

6

7 ¹Program in Computational Biology and Bioinformatics, Yale University

8 ²Department of Molecular Biophysics and Biochemistry, Yale University

9 ³Department of Computer Science, Yale University, 260/266 Whitney Avenue PO Box 208114,
10 New Haven, CT 06520, USA

11

12 * Correspondence should be addressed to M.G. (pi@gersteinlab.org)

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Abstract

Large-scale exome sequencing of tumors has enabled the identification of cancer drivers using recurrence and clustering-based approaches. Some of these methods also employ three-dimensional protein structures to identify mutational hotspots in cancer-associated genes. In determining such mutational clusters in structures, existing approaches overlook protein dynamics, despite the essential role of dynamics in protein functionality. In this work, we present a framework to identify driver genes using a dynamics-based search of mutational hotspot communities. After partitioning 3D structures into distinct communities of residues using anisotropic network models, we map variants onto the partitioned structures. We then search for signals of positive selection among these residue communities to identify putative drivers. We applied our method using the TCGA pan-cancer atlas missense mutation catalog. Overall, our analyses predict one or more mutational hotspots within the resolved structures of 434 genes. Ontological and pathway enrichment analyses implicate genes with predicted hotspots to be enriched in biological processes associated with tumor progression. Additionally, a comparison between our approach and existing hotspot detection methods that use structural data suggests that the inclusion of dynamics significantly increases the sensitivity of driver detection.

66 **Introduction**

67 Large-scale cancer genome studies such as The Cancer Genome Atlas (TCGA) project^{1,2} and the
68 International Cancer Genome Consortium (ICGC)^{3,4} have generated comprehensive catalogs of
69 somatic alterations for various cancer cohorts. The majority of these somatic variants incur little
70 or no functional consequence on tumor progression, and are thus often termed neutral
71 ‘passengers.’ In contrast, a handful of ‘driver’ mutations are considered to provide a selective
72 advantage to cancer cells. One of the critical goals of TCGA and ICGC projects has been to
73 distinguish between these positively selected “driver mutations”⁵⁻⁷ from a large number of
74 neutral passenger mutations.

75
76 A majority of the cancer driver detection algorithms quantify the recurrence of mutations to
77 identify significantly mutated genes and non-coding genomic elements⁸⁻¹¹. However, the somatic
78 mutation landscapes of cancer genomes are highly heterogeneous¹²⁻¹⁴ and exhibit a long tail of
79 low-frequency mutations^{11,13,15-17}. The presence of this long tail of rare somatic mutations, along
80 with limited cohort sizes, makes recurrence-based driver identification very challenging. An
81 alternative is to employ algorithms that aggregate mutation recurrence on gene/element-
82 levels^{18,19} or to predict the molecular functional impact of mutations²⁰ to distinguish drivers from
83 passengers. Compared to protein-truncating mutations and large structural variants, missense
84 mutations induce subtle changes, which are often difficult to interpret on the phenotypic level.
85 Thus, identifying missense driver mutations based on their molecular functional impact is also
86 challenging. In contrast, the signal of positive selection aggregated on functional elements or
87 sub-regions of the coding genome (such as protein domains²¹⁻²³, post-translational modification
88 sites (PTMS)²⁴⁻²⁶, protein interaction interfaces^{27,28} and mutation cluster/hotspots²⁹⁻³¹) has been
89 shown to be effective, despite their intrinsic limitations.

90
91 Prior studies have identified driver mutations based on their presence in mutational clusters²⁹⁻³¹,
92 which are sometimes called “hotspot” regions. These mutational clusters are defined based on
93 the proximity of somatic mutations within the primary sequence^{29,31} or three-dimensional
94 structure of a given protein³²⁻³⁶. Sequence-based mutation cluster identification algorithms^{29,31,37}
95 discover significantly mutated genes while considering an appropriate background mutation
96 model, trinucleotide context of mutations and distribution of silent mutations. However,

97 sequence-based approaches miss many hotspot regions, as they ignore spatial proximity between
98 residues that may be far apart in sequence but can be very close in 3-dimensional(3D) space^{38,39},
99 in the context of the fully-folded protein or protein ensembles. In contrast, despite being
100 inherently limited due to incomplete structural coverage of the proteome, 3D structure-based
101 mutational cluster definitions provide physical intuition or mechanistic insight into the roles of a
102 mutational cluster in cancer progression³²⁻³⁶. These structure-based methods compute residue
103 distances or generate residue-residue contact networks in the 3D structures of proteins to identify
104 a group of spatially proximal residues. Furthermore, mutation shuffling is performed to identify
105 significantly mutated residue clusters or hotspots on a protein structure. However, it is important
106 to note that, current approaches under this framework have failed to consider protein dynamics.

107
108 Proteins are inherently dynamic bio-molecules and sample large ensembles of conformations⁴⁰⁻
109 ⁴³. The energy landscape underlying the distribution of structures in these ensembles are often
110 altered based on external (thermodynamic)^{44,45} or internal (allosteric) signals^{43,46}. Previous
111 biophysical studies have clearly shown the crucial role of protein motions in conferring protein
112 functionality⁴⁷. Thus, one could argue that prior structure-based driver detection methods that
113 employ only the static structure of proteins are less sensitive when attempting to identify
114 functional residues through the mutation clustering approach.

115
116 In particular, a static crystalized structure provides only one limited snapshot of the protein, most
117 likely close to (or at) the bottom of the free energy landscape. In contrast, motion-weighted
118 community detection approach better reflects physical reality where proteins undergo two
119 general types of dynamics. First, a protein can dynamically oscillate around the bottom of the
120 energetic well or in second type of dynamics the underlying free energy landscape changes in
121 distinct ways, thereby shifting the protein conformation to an alternative functional state. In each
122 of these scenarios, communication between different communities plays a pivotal role in the
123 proper functioning of the protein. We posit that hotspot communities exist in large part because
124 certain select communities either play especially essential roles in these functional dynamics or
125 because their contributions to such dynamics are especially sensitive to mutations. Static
126 representations of protein structures presumably fail to define communities in light of their

127 essential roles in dynamics, and thus function. Furthermore, they potentially miss many critical
128 mutational clusters with a potential role in cancer progression.

129

130 In the current work, we address this issue by explicitly incorporating protein dynamics into our
131 new framework to identify mutational hotspot communities in protein structures. We applied this
132 framework to the TCGA pan-cancer atlas catalog of missense mutations to identify genes with
133 significantly mutated residue communities in protein structure. Our pan-cancer analysis
134 identifies 424 unique genes with at least one hotspot community in the corresponding protein
135 structure. The majority of these genes are involved in critical biological processes and pathways
136 involved in cancer progression including DNA repair, signal transduction, immune response,
137 apoptosis, and post-translational modifications. As expected, we observe higher cross-species
138 conservation score and greater functional impact scores for mutations present in these hotspot
139 communities. Furthermore, our prediction includes previously characterized driver genes with
140 hotspot communities in corresponding protein structure. Additionally, we also identify novel
141 genes with at least one hotspot community that were not detected by other mutation cluster
142 algorithms lacking protein dynamics information. Finally, we highlight some examples of driver
143 genes containing hotspot communities which are predicted to play a vital role in cancer
144 progression.

145

146 **Material and Methods**

147

148 **SNV dataset and mapping onto protein structure**

149 In this study, we leveraged the MC3(multiple-center mutation calling in multiple cancer)⁴⁸
150 somatic mutation dataset generated as part of the TCGA pancan atlas project. Briefly, the MC3
151 call set was generated using approximately 10,000 tumor/normal whole exome sequences
152 belonging to 33 different cancer types. Multiple callers, including MuTect⁴⁹, RADIA⁵⁰,
153 SomaticSniper⁵¹, and VarScan⁵² were applied to obtain high-confidence variant calls. Subsequent
154 filtering removed mutations due to lack of coverage, potential germline contamination, and other
155 artifacts. We utilized version 2.8 of the publicly accessible MC3 variant call set⁵. Furthermore,
156 we only analyzed missense mutations that were designated as ‘PASS’ based on the filtering
157 criterion. Moreover, we only analyzed variants from samples that were included in the whitelist

158 samples and were not hyper-mutated. This subset comprises 2.85 million mutations from 8937
159 samples in the pancan atlas project. Approximately 2.29 million mutations in this subset occupy
160 the coding regions of the genome that consists of 1.5 million missense mutations, 1.18 million
161 silent mutations, 0.6 million nonsense mutations, and 3.7K splice mutations.

162
163 We applied the Variant Annotation Tool (VAT)⁵³ to map TCGA missense mutations onto protein
164 structures. For each missense mutation, VAT provides an annotation that includes gene name,
165 transcript name, and the position of the residue getting affected in the translated protein
166 sequence. Additionally, it also provides the residue identity of the original and mutated residues.
167 Subsequently, we integrated VAT annotations with a BioMart⁵⁴ derived identifier map, which
168 consists of the gene identifier, transcript identifier, and the corresponding PDB ID, if available.
169 We restrict our analyses to mutations that map to crystal structures with resolution better than 3.0
170 Å. This restriction was applied to in order to most precisely identify residue communities in
171 protein structures. Overall, we mapped 0.329 million missense mutations on approximately
172 17,300 crystal structures in the current study.

173

174 **Workflow to identify three-dimensional hotspot communities in cancer**

175 As discussed above, our framework to predict driver genes through identification of hotspot
176 communities is novel compared to prior approaches as we explicitly include protein dynamics
177 information in our workflow (**Fig 1**). Briefly, our integrative workflow includes three distinct
178 components. First, we model large-scale conformational changes of a protein to identify dynamic
179 sub-regions of proteins (or “communities”). The large-scale conformational changes are modeled
180 using anisotropic network models (ANMs)^{46,55}. Subsequently, we model protein structure as a
181 residue-interaction network, where each residue constitutes a node in the network, and edges (or
182 connections between these nodes) form the physical interactions between these nodes.
183 Furthermore, edges in a network can be ‘weighted’ using the extent to which contacting residues
184 exhibit correlated movements within the dynamic structure of the protein. Highly correlated
185 motion (or movement vectors) between two residues that are physically in contact (though not
186 necessarily covalently linked) suggest that knowledge of the motions for one residue can provide
187 a great deal of information regarding the motions of the other residue. This mutual knowledge, in
188 a sense, suggests a strong degree of informational flow between residues. The weight for each

189 edge in the network corresponds to the “effective distance” of this edge, in which a strong degree
190 of correlated motion results in a short distance, and a weak correlation in the motions results in a
191 long distance. With this motion-weighted protein network, communities of residues are defined
192 with the Girvan-Newman algorithm⁵⁶. Communities are then defined as residue groups in which
193 each residue of a given community is connected to other residues of the community, and only
194 tangentially connected to residues outside the immediate community. These network-weighted
195 communities thus form densely inter-connected neighborhoods.

196
197 In order to identify mutational hotspot communities on a given protein structure, we mapped
198 missense mutations from TCGA cohorts onto three-dimensional protein structures.
199 Subsequently, we computed the frequency of mapped mutations for each community on the pan-
200 cancer level as well as in specific cancer cohorts. Furthermore, for each community with mapped
201 mutations, we performed a Fisher exact test to determine whether variants fall within a given
202 community is more frequently mutated than what would be expected by chance. This
203 significance test assigns an empirical p-value, which we correct for multiple hypothesis testing
204 using the Benjamini Hochberg method to identify significantly mutated hotspot communities on
205 protein structure for a given gene. We note that, for a substantial number of genes, there are
206 multiple PDB structures available. We remove this structural redundancy using structural
207 coverage (highest fraction of residues covered in the structure) as a filter to provide one to one
208 mapping between PDB structure and corresponding gene. The source code for the workflow is
209 available on the project’s Github page (<https://github.com/gersteinlab/HotComms>).

210

211

212 **Downstream Analyses**

213 We performed many downstream analyses to further validate our predictions. We extracted
214 PhyloP⁵⁷ and CADD⁵⁸ score for each mutation mapping onto protein structures. Furthermore, we
215 classified mutations into hotspot and non-hotspot mutations based on whether mutations are
216 mapped onto residues belonging to hotspot communities or otherwise. Subsequently, we
217 compared the phyloP score and CADD score distributions for hotspot and non-hotspot
218 mutations. We performed two-sided Kolmogorov-Smirnov(KS) test to assess the significance of
219 conservation score differences between hotspot and non-hotspot mutations. We apply the same

220 method to quantify such disparities for the molecular functional impact (CADD) score for
221 hotspot and non-hotspot mutations. Here, our null hypothesis is that the conservation or impact
222 score for hotspot and non-hotspot mutations are on average not different as they are being drawn
223 from the same distribution.

224
225 We also performed gene ontology(GO) enrichment and pathway enrichment analyses to further
226 validate the role of our putative driver genes in tumor progression. For the GO analysis, we
227 calculated the enrichment based on biological processes available from the GO database⁵⁹, and
228 we performed pathway enrichment analysis using the Reactome⁶⁰ as well as the KEGG
229 database⁶¹. We visualized the enrichment analysis result using the clusterProfiler⁶² package
230 available in Bioconductor.

231
232 Additionally, we also compared our predicted driver gene list derived from our hotspot
233 community analysis with other approaches that detect driver genes based on the presence of
234 mutation clusters on sequence or structure levels. One of the key differences between our
235 approach and other approaches is that we employ information on protein dynamics (along with
236 structural data) to determine hotspot communities. For structure-based methods, we obtained
237 driver gene list predicted from HotSpot3D³⁵, 3DHotSpot³⁴, HotMap³⁶ algorithms. All three of
238 these algorithms were previously applied on the TCGA Pancan Atlas data⁵, which allows us to
239 make meaningful comparisons with our work. However, we also note small differences in our
240 workflow compared to other structure-based approaches. For instance, HotMap tools employ
241 homology-model derived structures compared to other methods that rely only of experimentally
242 determined structure. Moreover, our method was applied only on crystal structure at higher
243 resolution compared to other methods that included NMR as well as crystal structures at higher
244 resolution. Finally, we also employed predicted driver genes from sequence-based cluster
245 analysis tool (OncodriverClust³¹) and previously curated driver genes in the cancer gene
246 census(CGC) database^{63,64}. We note that we excluded driver genes in CGC that play role in
247 cancer through INDELS, copy number aberrations or other structural variations. We used
248 UpsetR⁶⁵ package in R to visualized the multiway comparisons among predicted driver genes
249 from various tools and CGC database.

250

251 Finally, we also performed gene expression analysis to validate the role of our putative driver
252 genes in cancer at the transcriptome level. For this analysis, we obtained the TCGA RNA-Seq
253 quantification available for samples in the Pancan Atlas project². For each gene in our putative
254 driver gene list (based on hotspot community information), we compared the gene expression
255 distribution for sampled that harbored missense mutations to those that are not mutated. We
256 performed a two-sided KS test to evaluate the significance value for each gene in our putative
257 gene list. These significance tests were carried out separately for each cancer-type. However, we
258 combined the significance level(p-value) for each gene across multiple cancer types using the
259 Fisher method. We visualized significantly differentially expressed genes using a standard QQ
260 plot.

261

262 **Results**

263 **Pan-cancer analysis of genes containing mutations clusters**

264 We applied our workflow to identify significantly mutated hotspot communities for each cancer
265 cohort as well as on the pan-cancer level. As expected, we observed a comparatively higher
266 number of genes with at least one hotspot community on the pan-cancer level compared to
267 cancer-specific analysis. Our pan-cancer analysis identifies hotspot communities present on
268 protein structures of 434 unique genes (**Fig 2a, supplement table S1**). In contrast, a cancer-
269 specific analysis revealed 56 potential driver genes with 186 significantly mutated hotspot
270 community in the corresponding protein structure (**Supplement table S2**). Some of these genes
271 (including TP53, PIK3CA, BRAF, SPOP, KRAS, HRAS, and PTEN) have been previously
272 shown to be a driver for different cancer types. However, we also identified numerous novel
273 genes containing hotspot communities that might drive cancer progression. Previous studies
274 suggest that some of these novel genes including RHOC, NCOA1, and KLHL12 are involved in
275 various signaling pathways. Similarly, PSPC1, FOXO3, and XRCC5 are known to be pivotal for
276 immune response, apoptosis, and DNA repair, respectively. Furthermore, among these 434
277 genes, 12 genes had five or more hotspot communities whereas 352 genes had just one hotspot
278 community on their corresponding protein structure. These observations highlight the efficacy of
279 our approach in identifying novel and low-frequency putative driver genes with hotspot
280 communities.

281

282 Mutation cluster-based approaches assume that residues constituting such clusters are essential
283 for protein functions. Thus, a majority of cancer missense mutations occupying these hotspot
284 communities are very likely to disrupt the protein function. In order to validate this assumption,
285 we quantified the cross-species conservation measure (PhyloP score⁵⁷) for mutations in hotspot
286 as well as non-hotspot communities on protein structures. As expected, we observe higher
287 average conservation score for mutations mapping to residues in hotspot communities compared
288 to those, which are present outside. Furthermore, the observed difference in conservation is
289 statically significant (two-sided KS test, p-value < 2e-5) (**Fig 2b**). Similarly, the putative
290 molecular functional impact (CADD score⁵⁸) of mutations occupying hotspot communities was
291 significantly higher compared to those mapping to non-hotspot communities (two-sided KS test,
292 p-value < 2e-5) (**Fig 2c**).

293
294 We also preformed gene ontology⁶² and pathway enrichment analysis to decipher the biological
295 function of genes with predicted hotspot communities. The biological process based gene
296 ontology enrichment analysis indicate role of putative driver genes in diverse biological function
297 including immune response, cell differentiation, kinase activities, post-translational
298 modifications, apoptosis and DNA repair (**Fig 2d & Supplement table S3**). Similarly, reactome
299 pathway⁶⁰ based enrichment analysis suggest role of putative driver genes with hotspot
300 communities in various signaling pathways (**Supplement table S4**) including NTRK signaling,
301 DAP12 signaling, EGFR signaling and MAP kinase-associated signaling. Additionally, these
302 genes are also enriched among DNA repair and non-homologous end-joining associated
303 pathways (**Fig 2e**). Furthermore, KEGG pathway⁶⁶ based enrichment analysis indicate role of our
304 putative driver genes in various cancer subtypes (bladder, pancreatic, breast, CML, melanoma,
305 AML, glioma) (**Supplement Fig1 & Supplement table S5**).

306

307 **Comparison of 3D structure based clustering methods**

308 We performed consensus analysis between our approach to the driver genes curated in the
309 COSMIC⁶⁷ database. Furthermore, we also performed a comparison between putative driver
310 genes identified using our workflow and genes identified as drivers by other mutation cluster
311 detection algorithms that do not take protein dynamics into account. The majority of these
312 additional algorithms employ the three-dimensional structure of a protein to identify mutational

313 cluster except the OncoDriveClust³¹ tool, which searches for hotspot mutations on the sequence
314 level. Overall, our workflow identified many additional genes (288 genes) with hotspot
315 communities compared to other mutation hotspot analysis tools (**Fig 3a**). One exception being
316 the HOTMAP³⁶ algorithm that utilizes protein homology model in addition to protein structure.
317 Thus, it identifies significantly higher number of unique genes (620 genes) with mutation cluster
318 compared to any other tool. Furthermore, our approach identified 146 genes (34% of our gene
319 list) with hotspot communities that are either curated as a driver gene in COSMIC or predicted to
320 contain a mutation cluster by another tool (**Fig 3a**). Among these 146 genes, 89 genes
321 overlapped with putative driver genes identified by HOTMAP algorithm, whereas 63 genes
322 overlapped with drivers in COSMIC. As expected, we observed the lowest overlap (33 genes,
323 7% of our putative driver gene list) with sequence-based method (OncoDriveClust; **Fig 3a**).

324
325 Additionally, we analyzed TCGA expression data to obtain additional evidence corroborating the
326 biological validity of putative driver genes identified through our workflow. Intuitively, one
327 would expect a significant difference in gene expression level between samples with and without
328 mutation for genes that were predicted to contain a significantly mutated hotspot community. For
329 each candidate gene, we quantified the statistical significance in expression distribution
330 differences using two-sided KS test. Furthermore, we performed this test for individual cancer
331 type, and the corresponding p-values were combined across cancer types using Fisher's method
332 to provide a pan-cancer significance measure. Overall, our analysis identified 60 genes including
333 TP53(p-value 3.59e-66), SPTA1 (p-value 8.58e-32), PIK3CA (p-value 7.06e-25), KRAS (p-
334 value 5.73e-11), and EGFR (p-value 2.78e-06) that were differentially expressed across cancer
335 types (**Fig 3b & Supplement table S6**). A subset of these differentially expressed genes such as
336 MYH7 (p-value 4.22e-15), ROS1 (p-value 3.26e-13), TIAM1 (p-value 2.48e-12), PTPRD (p-
337 value 3.96e-23), and HUWE1 (p-value 4.84e-10) are potentially novel driver genes with
338 predicted hotspot communities (**Fig 3b & Supplement table S6**). Moreover, we note that 76%
339 of our putative driver gene list with significantly mutated hotspot communities were
340 differentially expressed in at least one TCGA cancer cohort.

341
342 Finally, we performed GO and pathway enrichment analysis on novel genes that we predict to
343 contain mutational hotspot communities. However, these genes were neither present in the

344 COSMIC driver database nor were predicted to encompass mutation cluster through other
345 hotspot identification tools. We observed significant enrichment of these genes in crucial
346 biological processes (**Supplement table S7**) including DNA conformation change, regulation of
347 immune response, regulation of stem cell differentiation, nucleosome organization, and
348 endothelial cell apoptotic process (**Supplement Fig2**). Similarly, pathway enrichment analysis
349 implicates their role in DNA repair, SUMOylation, RHO GTPase activity, telomere
350 maintenance, and various signaling pathways (**Fig 3c & Supplement table S8**).

351

352 **Case studies highlighting the roles of hotspot communities in deciphering driver** 353 **mechanisms**

354 Integration of protein 3D-structure and protein dynamics to identify driver genes has a clear
355 advantage over other methods that do not leverage protein structure or protein dynamics
356 information. Our method allows us to investigate disruption in protein structure and function
357 induced by missense mutations that occupy within predicted hotspot communities. We also note
358 that the majority of our hotspot communities encompass residues that are pivotal for important
359 protein functions including allostery, bimolecular signaling, protein binding, and post-translation
360 modifications. The sensitive detection of functional sites on protein structure helps to decipher
361 the underlying biophysical mechanism that plays a crucial role in cancer growth. Here, we
362 highlight three examples testifying the utility of our framework in gaining biophysical insight
363 into cancer progression through disruption of predicted hotspot communities. These examples
364 include an oncogene(BRAF), tumor suppressor gene(PIK3R1), and a novel putative driver
365 gene(PTPRD) that are predicted to contain multiple hotspot communities on their respective
366 protein structure.

367

368 **Missense hot spot communities: PIK3R1**

369 The PI3KR1 gene encodes the alpha subunit of the enzyme Phosphatidylinositol 3-kinase
370 regulatory, which plays a crucial role in a variety of cellular processes including cell survival,
371 regulation of gene expression, cell metabolism and cytoskeletal rearrangement⁶⁸. Mutations in
372 PIK3KR1 gene has previously been implicated as a tumor suppressor gene in breast cancer.
373 Recent therapeutic studies have targeted PI3K inhibition resulting in a decrease in cellular
374 proliferation and reduced metastasis in the mouse model. PI3Ks are obligate heterodimers

375 composed of a p110 subunit and a regulatory subunit. Previous studies have identified four
376 distinct domains belonging to the catalytic P110 alpha subunit that harbor somatic mutations
377 leading to an increase in PI3K activity. We observe two distinct hotspot communities (**Fig 4a**) on
378 the co-crystal structure (PDB ID: 2V1Y) of the protein complex that comprises ABD domain
379 of the P110 alpha subunit and the iSH2 domain of the p85 alpha regulatory subunit. The two
380 hotspot communities are composed of 28(community 5) and 26(community 7) residues,
381 respectively (Fig 4a). On the pan-cancer level, we observe 24 and 16 mutations that map to
382 community 5 and community7 on the co-crystal structure, respectively. These distinct hotspot
383 communities are adjacent to each other in the same helical structure. However, we observe a
384 small kink in this helical structure, which presumably lead to distinct protein motions associated
385 with these two different hotspot communities.

386

387 *Missense hotspot communities in BRAF gene*

388 BRAF gene encodes a protein belonging to the serine/threonine protein kinase family that
389 regulates MAP kinase and ERK signaling pathway⁶⁹. This pathway is considered to be essential
390 for a number of biological functions including cell differentiation, cellular growth, senescence,
391 and apoptosis. Somatic mutations in the BRAF gene are often implicated in various cancer
392 subtypes including melanoma, colorectal cancer, prostate cancer, non-small-cell lung cancer, and
393 papillary thyroid tumors. It has been proposed that BRAF induce dysregulation in the binding of
394 Ras proteins to Raf and MEK proteins in the Ras/RAF/MEK/ERK signaling cascade that leads to
395 over-activation of the signaling pathway and subsequent oncogenesis. Multiple enzyme
396 inhibitors have been designed to target BRAF kinase in the tumor. One such inhibitor SB-
397 590885 has been co-crystallized with BRafV600E kinase domain at the X-ray resolution of 2.9
398 Angstrom (PDBID: 2FB8)⁷⁰. A previous study indicates the role of pi-stacking interactions,
399 hydrogen bonds and salt bridges in stabilizing the interaction between these two subunits in the
400 crystal structure. In our study, we identified one hotspot community in this co-crystal structure
401 (**Fig 4b**). This hotspot community is composed of 52 residues that constitute a beta sheet
402 secondary structure. Interestingly, we also observe that SB-590885 inhibitor occupies the same
403 hotspot community.

404

405 *Missense hotspot community in TPRD gene*

406 The PTPRD gene encodes a protein that belongs to the protein tyrosine phosphatase(PTP)
407 family. PTP proteins are considered essential for regulating cellular proliferation, differentiation,
408 and oncogenic transformation. PTPRD gene encodes a transmembrane protein containing a
409 cytoplasmic tyrosine phosphatase domain. Previous studies have shown that PTPRD genes are
410 frequently deleted in various cancer types including glioma, neuroblastoma, and lung cancer⁷¹.
411 However, we note PTPRD is not identified as missense driver in cosmic catalog. Moreover,
412 previous studies did not identify presence of mutational hotspot communities in the PTPRD
413 gene. In contrast, our analysis identifies one hotspot community in the crystal structure (PDB ID:
414 2YD7) of the receptor protein tyrosine phosphatase(RPTP) sigma subunit. RPTPs are cell
415 surface proteins with intracellular PTP activity and extracellular domains that are sequentially
416 homologous to cell adhesion molecules. Moreover, RPTP sigma subunit is considered necessary
417 for nervous system development and function. In our analysis, somatic mutations mapped to two
418 communities (community 2 & 4) on the crystal structure of the RPTP sigma subunit. Our
419 workflow predicts one hotspot community that comprise of 47 residues in the crystal structure of
420 PTPRD (**Fig 4c**) and adopts a beta strand conformation.

421

422

423 **Discussion**

424 The underlying heterogeneous characteristic⁷² of cancer makes interpretability of genomic
425 alterations in a cancer genome very challenging. In particular, genomic heterogeneity poses a
426 major challenge in identifying key driver mutations in cancer. Large-scale cancer genome
427 sequencing efforts have helped us to generate comprehensive catalogs of driver mutations⁵ in
428 various cancer types. However, the canonical recurrence-based driver detection algorithms have
429 failed to identify low-frequency or rare drivers. The limited cohort size¹¹ and heterogeneity¹⁴ in
430 cancer genome provides limited power to identify low-frequency drivers using the canonical
431 position level recurrence algorithms. A simplistic approach to address the issue of missing rare
432 driver will be to sequence more patients for a given cancer type. However, this approach will be
433 particularly challenging for highly heterogeneous cancer cohorts with multiple subtypes⁷³ within
434 a cancer type. Moreover, this approach will not be practical for certain rare cancers including
435 neuroblastoma, angiosarcoma, Hodgkin's lymphoma, and others. A suitable alternative is to
436 quantify recurrence over functional elements or sub-gene levels⁷⁴ such as post-translational

437 modification sites (PTMS)^{25,26}, protein interaction interfaces²⁸ and mutational clusters^{33–36,38}. In
438 particular, many driver detection algorithms search for the presence of mutational hotspot on the
439 3D-protein structures to identify putative driver genes. Compared to sequence-based driver
440 detection methods, using protein structural data can help to decipher the underlying molecular
441 mechanisms that influence cancer progression. However, current approaches to identify cancer
442 mutation hotspots on protein structure and corresponding driver genes completely ignore the role
443 of protein dynamics, which is considered essential for protein function. Thus, here we propose a
444 new framework that utilizes protein dynamics along with the 3D-structure of proteins to identify
445 missense hotspot communities on protein structure and corresponding putative driver genes.

446
447 Overall, our workflow identified 802 hotspot communities on crystal structures of proteins
448 corresponding to 434 unique genes on the pan-cancer level. We also compared our putative
449 driver gene list with previous experimental and prediction studies derived driver gene list.
450 Among our putative driver gene list, we find 36% of genes are either known or predicted to be
451 driver genes based on previous studies. We term the remaining 64% of genes as novel drivers in
452 our study. We performed many downstream analyses on our putative driver genes to highlight
453 their role in cancer progression. Our framework assumes that a residue community on a protein
454 structure represents a putative functional subunit of a protein. Thus, high mutation densities in
455 such communities (compared to a random expectation) is very likely to alter protein function.
456 One would expect that mutations influencing residues in these communities will have a high
457 functional impact as they can drive cancer progression. Our observation is consistent with this
458 hypothesis, as we find that missense mutations occupying hotspot communities in proteins
459 structures are highly conserved across species and have a higher molecular functional impact
460 compared to those outside such hotspot communities.

461
462 Furthermore, we also observe significantly high enrichment of our putative driver genes with
463 predicted hotspot communities in vital biological processes and pathways that are relevant for
464 oncogenesis. For instance, ontology analysis indicates enrichment of our putative driver genes in
465 biological processes associated with regulation and activation of innate immune response. This
466 observation is consistent with the current notion that dysfunction in immune response
467 contributed through genomic alterations will allow tumor cells to evade immune detection due to

468 lack of effective immune response. Additionally, we also observe a significant enrichment of
469 putative driver genes in cell differentiation and cell growth processes, such as the regulation of
470 hematopoiesis and myeloid cell differentiation, which were previously implicated in tumor
471 growth. Moreover, we observed a high enrichment of our putative driver genes in the regulation
472 of kinase activities including protein serine/threonine and MAP kinase activities. Additionally,
473 these genes are also enriched among ERK1/ERK2 signaling cascade, protein kinase B signaling,
474 PI3K/AKT signaling, FGFR1 signaling, NTRK1 signaling, apoptosis signaling, and various
475 other signaling pathways. Presence of aberrant signaling pathways is an essential hallmark of
476 cancer. Thus, enrichment of our putative genes in critical signaling pathways provides clear
477 biological evidence for their role in cancer. Moreover, these genes are enriched for DNA repair
478 function via non-homologous end joining(NHEJ) and other non-recombination based repair
479 mechanisms. Finally, we note that we observed the same enrichment for the subset of novel
480 genes in our putative driver gene list, that have not been identified as drivers in previous studies.

481
482 Genomic alterations that are consequential for tumor growth are often manifested on the
483 transcriptome level such that mutated driver genes are often differentially expressed compared to
484 a healthy population or patients without any mutation in driver genes. We leveraged the
485 transcriptome data from TCGA to further validate out predicted driver genes based on hotspot
486 community identification. We identified 60 genes among our predicted driver genes that were
487 significantly differentially expressed in tumor samples with missense mutations in those genes
488 compared to those without among multiple cancer cohorts. These differentially expressed driver
489 genes include novel as well previously established driver genes. Similar to genetic data,
490 transcriptomic data in TCGA is limited for specific cancer cohort that provides insufficient
491 power to identify all differentially expressed genes. However, we note that 76% of our putative
492 driver genes were differentially expressed in at least one TCGA cancer cohort. These analyses
493 further validate our hotspot community-based driver detection approach.

494
495 In the context of investigating the molecular mechanism underlying tumor growth, protein
496 structure-based driver detection methods offer significant advantages over approaches that are
497 only sequence-based. However, structure-based methods suffer from limited coverage of the
498 human proteome. Thus, the applicability of structure-based methods is inherently limited only to

499 mutations that can be mapped onto protein structure. A prior study³⁶ has applied homology
500 model derived structures to circumvent the issue of limited structural coverage. However, the
501 accuracy of homology-based models has shown to be limited for various protein complexes and
502 transmembrane proteins. Moreover, modeling protein motions for homology-model derived
503 proteins structures will be most likely less accurate thus affecting the sensitivity of our approach.
504 Nevertheless, significant technical improvement in crystallographic and cryoEM techniques⁷⁵ are
505 expected to expand the current structurally-resolved proteome. In particular, cryoEM
506 technologies⁷⁵ now allows us to obtain a high-resolution structure of large-size proteins and other
507 biomolecular complexes that were previously elusive. Thus, we anticipate an essential role of our
508 approach in future studies aimed at discovering low-frequency drivers in various cancer cohorts.
509 Additionally, knowledge of protein motions (along with structure) can potentially help in
510 uncovering drug interaction with hotspot communities. Such studies are likely to open new
511 therapeutic avenues for various cancers and will help us realize the goal of precision medicine in
512 cancer.

513

514

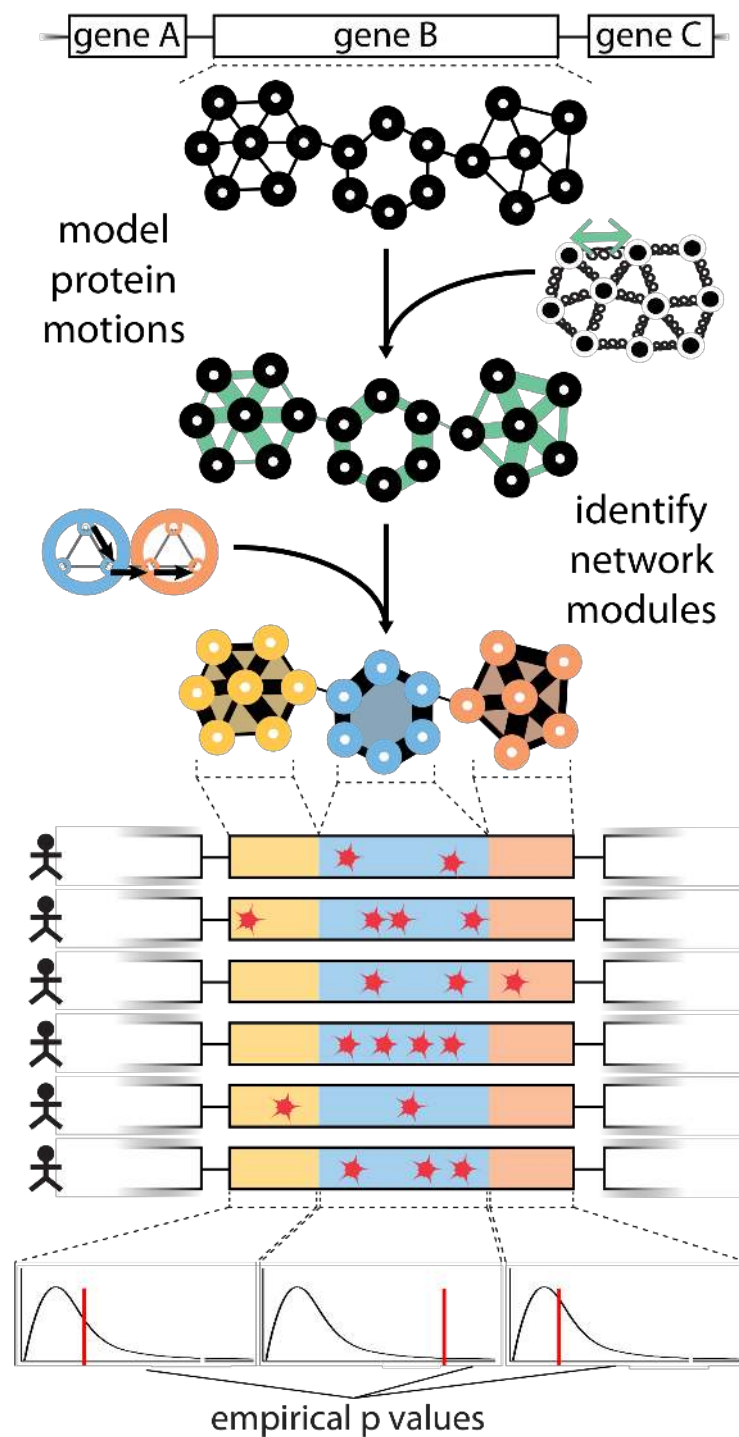
515 **Acknowledgments**

516 We acknowledge support from the NIH and from the AL Williams Professorship funds. We also
517 acknowledge help of Jonathan Warrell and Timur Galeev for providing valuable feedbacks for
518 improving the manuscript.

519

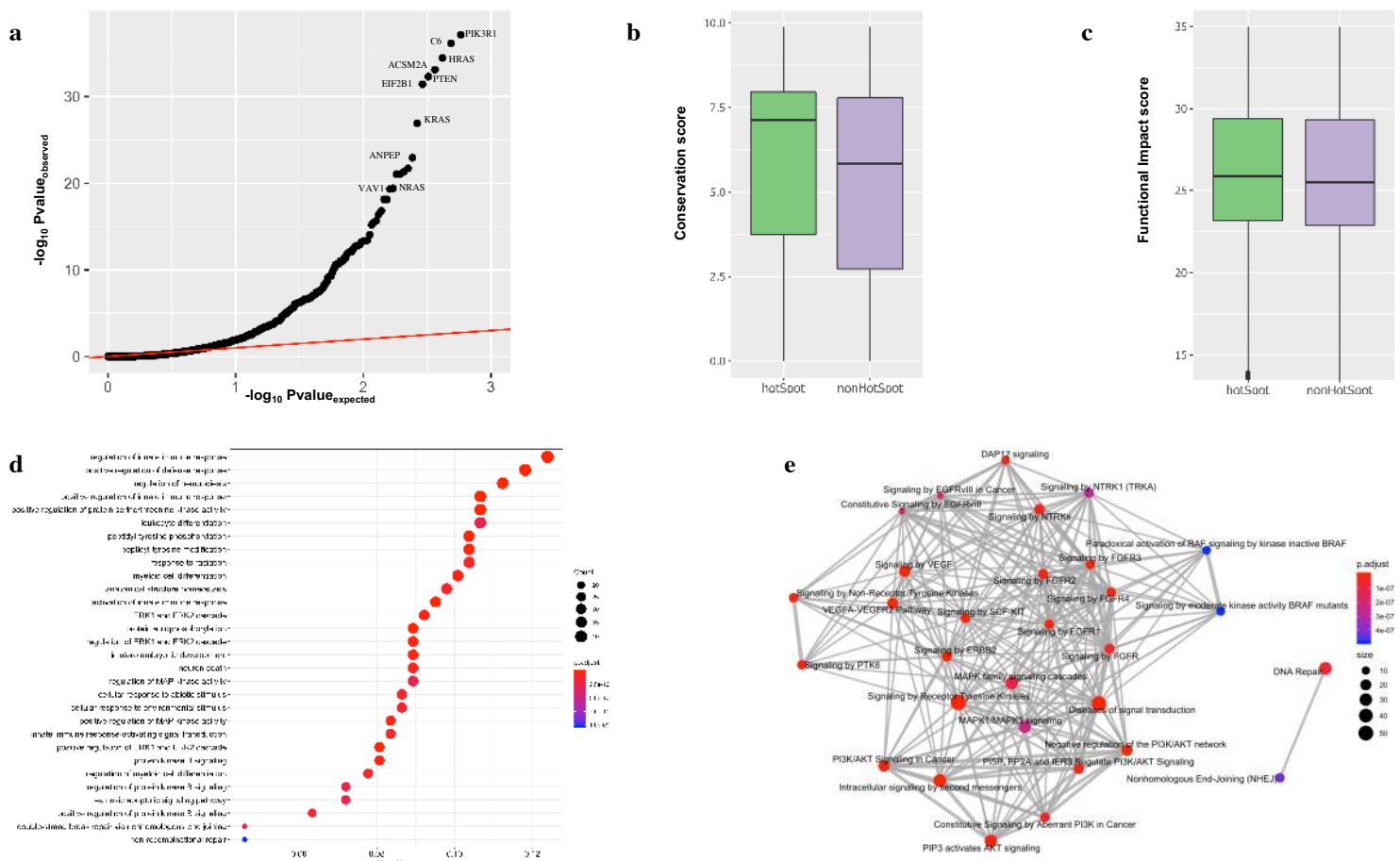
520

521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547



548 **Fig 1. Workflow of HotCommics to identify putative driver genes:** This integrative approach
549 utilizes protein community information along with mapped mutations onto protein structure to
550 identify significantly mutated communities in protein structure. Fisher method is employed to
551 quantify significance value for each community with mapped mutations.

552



553

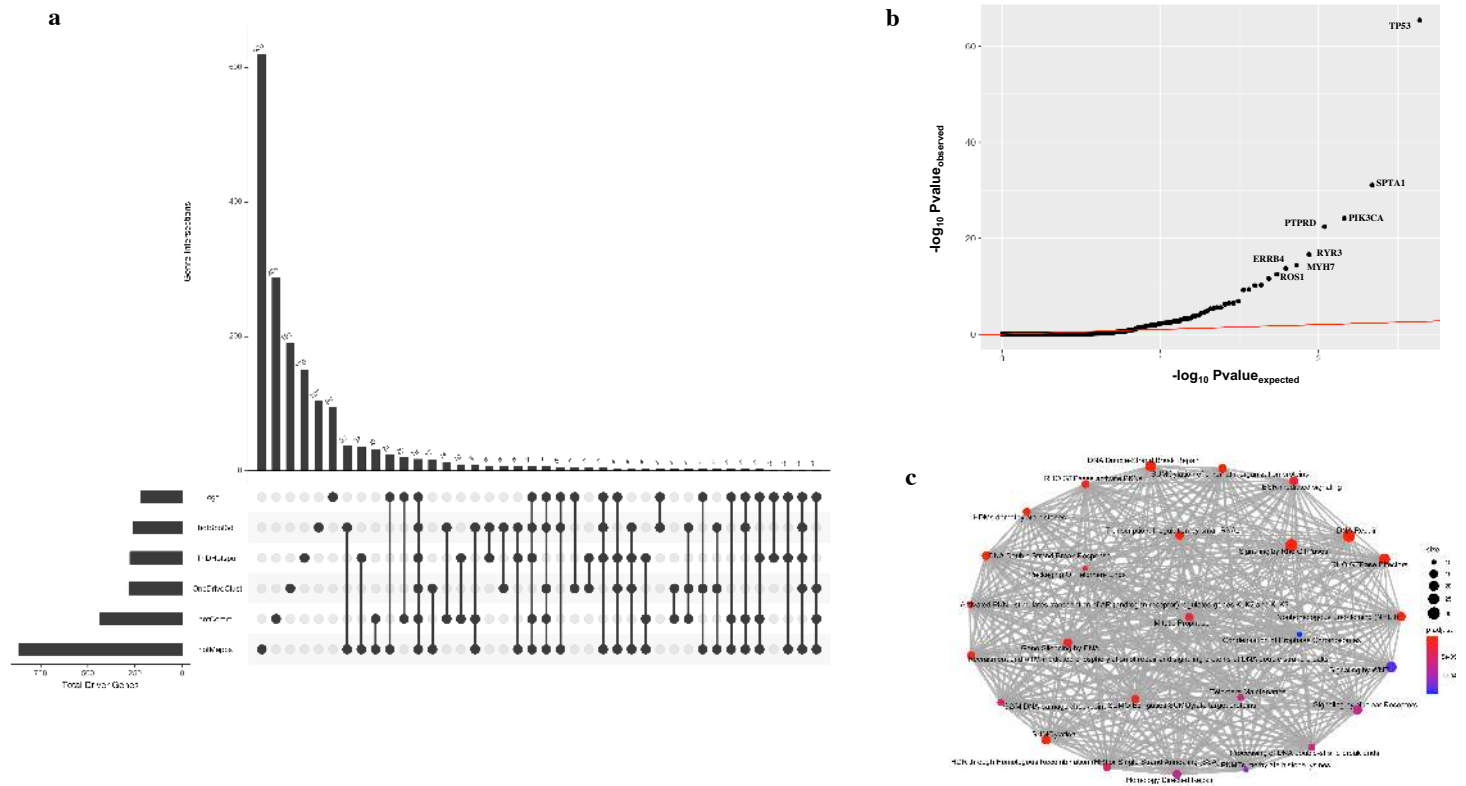
554

555

556 **Fig2. Pan-cancer analysis of putative driver genes with hotspot communities:** a) pan-cancer
 557 q-q plot for genes with hotspot communities, b) PhyloP conservation score comparison between
 558 mutations occupying hotspot communities against non-hotspot communities on protein
 559 structures, c) CADD score correlation between mutations occupying hotspot communities
 560 against non-hotspot communities on protein structures, d) Biological process enrichment analysis
 561 for putative driver genes with at least one hotspot. X-axis corresponds to gene ratio that
 562 corresponds to the fraction of putative driver genes belonging to a particular biological process.
 563 The color code and size correspond to corrected p-value and number of genes involved in the
 564 biological process, respectively, e) Reactome based pathway enrichment analysis. The color code
 565 and size correspond to corrected p-value and number of genes involved in the biological process,
 566 respectively.

567

568



569

570

571

572 **Fig3. Pan-cancer analysis of putative driver genes with hotspot communities:** a) Comparison
 573 of multiple driver detection algorithms including HotCommics. We used the most recent version
 574 of the Cancer Gene Census database for this analysis. Remaining algorithms were also run on the
 575 MC3 variant call set, b) Q-q plot highlighting differentially expressed putative driver genes
 576 across multiple cancer types, c) Pathway level enrichment analysis of singleton genes identified
 577 by the HotCommics algorithm that was novel for putative driver genes identified by other
 578 algorithms and CGC database.

579

580

581

582

583

584

585

586

587

588

a

589

590

591

592

593

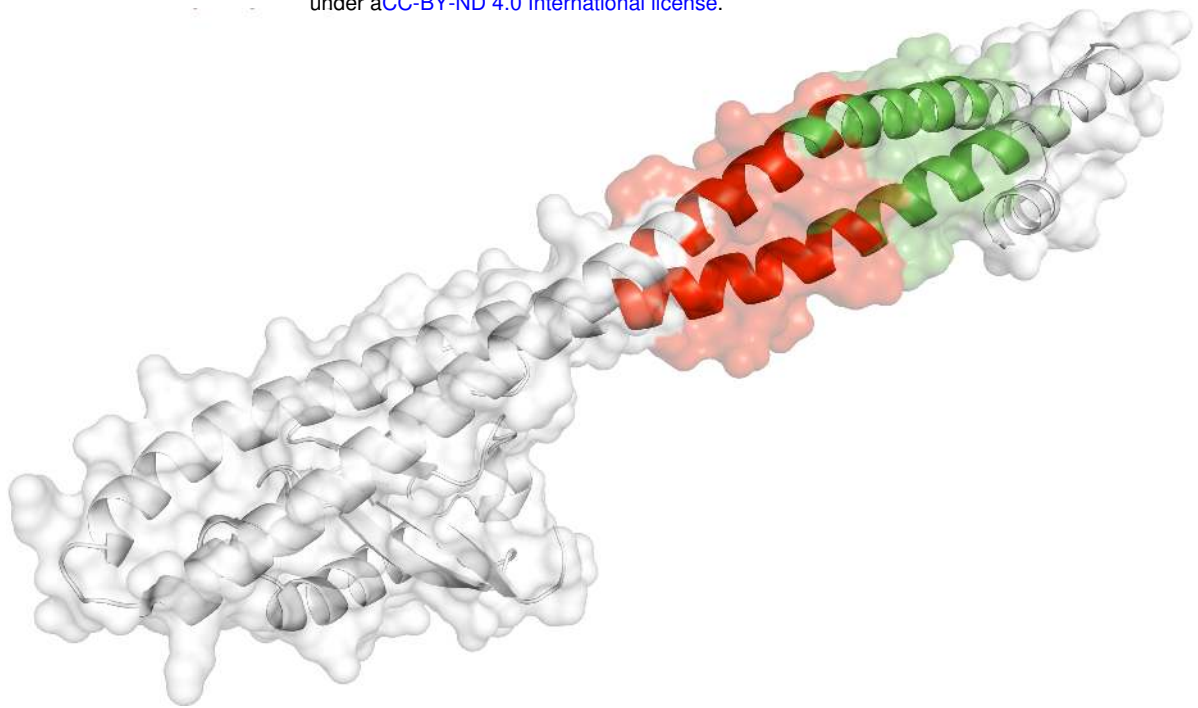
594

595

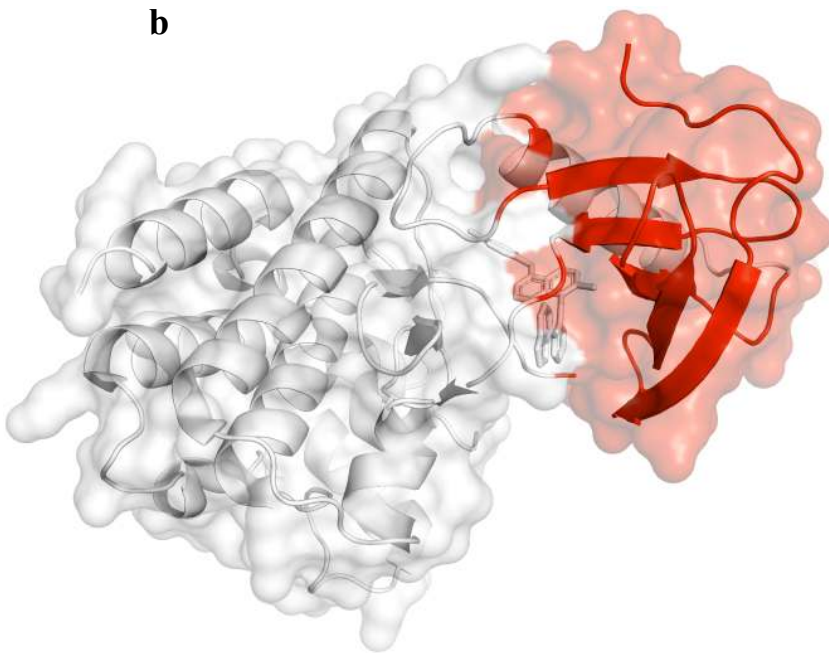
596

597

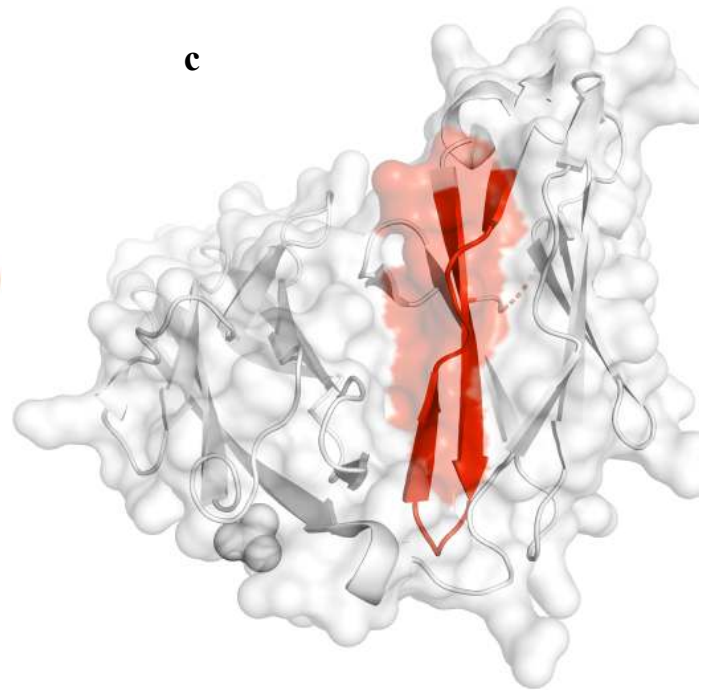
598



b



c



599

600 **Fig4. Examples of TSG, oncogene, and putative driver genes with hotspot communities:**

601 a) Example hotspot communities (shown in red) on the PIK3R1 gene as identified by our
602 workflow. We note that previous studies have identified the PIK3R1 gene as a tumor suppressor

603 gene, b) Example hotspot communities (shown in red) on the BRAF gene as identified by our
604 workflow. We note that previous studies have identified BRAF1 gene as an oncogene, c)

605 Example hotspot communities (shown in red) on the PTPRD gene as identified by our workflow.

606 We note PTPRD is an example of novel putative driver genes with hotspot community with
607 significant differential gene expression.

608 **References**

- 609 1. Weinstein, J. N. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.*
610 **45**, 1113–20 (2013).
- 611 2. Ding, L. *et al.* Perspective on Oncogenic Processes at the End of the Beginning of Cancer
612 Genomics. *Cell* **173**, 305–320.e10 (2018).
- 613 3. International Cancer Genome Consortium *et al.* International network of cancer genome
614 projects. *Nature* **464**, 993–8 (2010).
- 615 4. Campbell, P. J. *et al.* Pan-cancer analysis of whole genomes. *bioRxiv* 162784 (2017).
616 doi:10.1101/162784
- 617 5. Matthew Bailey, A. H. *et al.* Comprehensive Characterization of Cancer Driver Genes and
618 Mutations Article Comprehensive Characterization of Cancer Driver Genes and
619 Mutations. *Cell* **173**, 371–376.e18 (2018).
- 620 6. Rheinbay, E. *et al.* Discovery and characterization of coding and non-coding driver
621 mutations in more than 2,500 whole cancer genomes. *bioRxiv* 237313 (2017).
622 doi:10.1101/237313
- 623 7. Sabarinathan, R. *et al.* The whole-genome panorama of cancer drivers. *bioRxiv* 190330
624 (2017). doi:10.1101/190330
- 625 8. Ding, L., Wendl, M. C., McMichael, J. F. & Raphael, B. J. Expanding the computational
626 toolbox for mining cancer genomes. *Nat. Rev. Genet.* **15**, 556–570 (2014).
- 627 9. Raphael, B. J., Dobson, J. R., Oesper, L. & Vandin, F. Identifying driver mutations in
628 sequenced cancer genomes: Computational approaches to enable precision medicine.
629 *Genome Medicine* **6**, (2014).
- 630 10. Tamborero, D. *et al.* Comprehensive identification of mutational cancer driver genes
631 across 12 tumor types. *Sci. Rep.* **3**, 2650 (2013).
- 632 11. Lawrence, M. S. *et al.* Discovery and saturation analysis of cancer genes across 21 tumour
633 types. *Nature* **505**, 495–501 (2014).
- 634 12. Garraway, L. A. & Lander, E. S. Lessons from the cancer genome. *Cell* **153**, 17–37
635 (2013).
- 636 13. Stratton, M. R. Exploring the genomes of cancer cells: Progress and promise. *Science* **331**,
637 1553–1558 (2011).
- 638 14. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-
639 associated genes. *Nature* **499**, 214–218 (2013).
- 640 15. Armenia, J. *et al.* The long tail of oncogenic drivers in prostate cancer. *Nat. Genet.* **50**,
641 645–651 (2018).
- 642 16. Greenman, C. *et al.* Patterns of somatic mutation in human cancer genomes. *Nature* **446**,
643 153–158 (2007).
- 644 17. Beerwinkler, N. *et al.* Genetic progression and the waiting time to cancer. *PLoS Comput.*
645 *Biol.* **3**, 2239–2246 (2007).
- 646 18. Lawrence, M. S. *et al.* Discovery and saturation analysis of cancer genes across 21 tumour
647 types. *Nature* **505**, 495–501 (2014).
- 648 19. Dees, N. D. *et al.* MuSiC: Identifying mutational significance in cancer genomes. *Genome*
649 *Res.* **22**, 1589–1598 (2012).
- 650 20. Gonzalez-Perez, A. & Lopez-Bigas, N. Functional impact bias reveals cancer drivers.
651 *Nucleic Acids Res.* **40**, (2012).
- 652 21. Nehrt, N. L., Peterson, T. A., Park, D. H. & Kann, M. G. Domain landscapes of somatic
653 mutations in cancer. *BMC Genomics* **13 Suppl 4**, (2012).

- 654 22. Peterson, T. A., Gauran, I. I. M., Park, J., Park, D. H. & Kann, M. G. Oncodomains: A
655 protein domain-centric framework for analyzing rare variants in tumor samples. *PLoS*
656 *Comput. Biol.* **13**, (2017).
- 657 23. Yang, F. *et al.* Protein Domain-Level Landscape of Cancer-Type-Specific Somatic
658 Mutations. *PLoS Comput. Biol.* **11**, (2015).
- 659 24. Reimand, J. & Bader, G. D. Systematic analysis of somatic mutations in phosphorylation
660 signaling predicts novel cancer drivers. *Mol. Syst. Biol.* **9**, (2013).
- 661 25. Narayan, S., Bader, G. D. & Reimand, J. Frequent mutations in acetylation and
662 ubiquitination sites suggest novel driver mechanisms of cancer. *Genome Med.* **8**, (2016).
- 663 26. Reimand, J., Wagih, O. & Bader, G. D. The mutational landscape of phosphorylation
664 signaling in cancer. *Sci. Rep.* **3**, (2013).
- 665 27. Porta-Pardo, E. & Godzik, A. E-Driver: A novel method to identify protein regions
666 driving cancer. *Bioinformatics* **30**, 3109–3114 (2014).
- 667 28. Porta-Pardo, E., Garcia-Alonso, L., Hrabe, T., Dopazo, J. & Godzik, A. A Pan-Cancer
668 Catalogue of Cancer Driver Protein Interaction Interfaces. *PLoS Comput. Biol.* **11**, (2015).
- 669 29. Miller, M. L. *et al.* Pan-Cancer Analysis of Mutation Hotspots in Protein Domains. *Cell*
670 *Syst.* **1**, 197–209 (2015).
- 671 30. Van den Eynden, J., Fierro, A. C., Verbeke, L. P. C. & Marchal, K. SomInaClust:
672 Detection of cancer genes based on somatic mutation patterns of inactivation and
673 clustering. *BMC Bioinformatics* **16**, (2015).
- 674 31. Tamborero, D., Gonzalez-Perez, A. & Lopez-Bigas, N. OncodriveCLUST: Exploiting the
675 positional clustering of somatic mutations to identify cancer genes. *Bioinformatics* **29**,
676 2238–2244 (2013).
- 677 32. Ryslik, G. A. *et al.* A spatial simulation approach to account for protein structure when
678 identifying non-random somatic mutations. *BMC Bioinformatics* **15**, (2014).
- 679 33. Kamburov, A. *et al.* Comprehensive assessment of cancer missense mutation clustering in
680 protein structures. *Proc. Natl. Acad. Sci. U. S. A.* **112**, E5486-95 (2015).
- 681 34. Gao, J. *et al.* 3D clusters of somatic mutations in cancer reveal numerous rare mutations as
682 functional targets. *Genome Med.* **9**, (2017).
- 683 35. Niu, B. *et al.* Protein-structure-guided discovery of functional mutations across 19 cancer
684 types. *Nat. Genet.* **48**, 827–837 (2016).
- 685 36. Tokheim, C. *et al.* Exome-scale discovery of hotspot mutation regions in human cancer
686 using 3D protein structure. *Cancer Res.* **76**, 3719–3731 (2016).
- 687 37. Ye, J., Pavlicek, A., Lunney, E. A., Rejto, P. A. & Teng, C. H. Statistical method on
688 nonrandom clustering with application to somatic mutations in cancer. *BMC*
689 *Bioinformatics* **11**, (2010).
- 690 38. Ryslik, G. A., Cheng, Y., Cheung, K. H., Modis, Y. & Zhao, H. A graph theoretic
691 approach to utilizing protein structure to identify non-random somatic mutations. *BMC*
692 *Bioinformatics* **15**, (2014).
- 693 39. Ryslik, G. A., Cheng, Y., Modis, Y. & Zhao, H. Leveraging protein quaternary structure
694 to identify oncogenic driver mutations. *BMC Bioinformatics* **17**, (2016).
- 695 40. Frauenfelder, H., Sligar, S. & Wolynes, P. The energy landscapes and motions of proteins.
696 *Science (80-.)*. **254**, 1598–1603 (1991).
- 697 41. Tsai, C.-J. & Nussinov, R. The free energy landscape in translational science: how can
698 somatic mutations result in constitutive oncogenic activation? *Phys. Chem. Chem. Phys.*
699 **16**, 6332–41 (2014).

- 700 42. Boehr, D. D., Nussinov, R. & Wright, P. E. The role of dynamic conformational
701 ensembles in biomolecular recognition. *Nature Chemical Biology* **5**, 789–796 (2009).
- 702 43. Nussinov, R. & Tsai, C.-J. Allosteric in disease and in drug discovery. *Cell* **153**, 293–305
703 (2013).
- 704 44. Onuchic, J. N., Luthey-Schulten, Z. & Wolynes, P. G. Theory of protein folding: the
705 energy landscape perspective. *Annu. Rev. Phys. Chem.* **48**, 545–600 (1997).
- 706 45. Tsai, C.-J. & Nussinov, R. The free energy landscape in translational science: how can
707 somatic mutations result in constitutive oncogenic activation? *PCCP* 6332–41 (2014).
- 708 46. Clarke, D. *et al.* Identifying Allosteric Hotspots with Dynamics: Application to Inter- and
709 Intra-species Conservation. *Structure* (2016). doi:10.1016/j.str.2016.03.008
- 710 47. Henzler-Wildman, K. & Kern, D. Dynamic personalities of proteins. *Nature* **450**, 964–972
711 (2007).
- 712 48. Ellrott, K. *et al.* Scalable Open Science Approach for Mutation Calling of Tumor Exomes
713 Using Multiple Genomic Pipelines. *Cell Syst.* **6**, 271–281.e7 (2018).
- 714 49. Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and
715 heterogeneous cancer samples. *Nat. Biotechnol.* **31**, 213–219 (2013).
- 716 50. Radenbaugh, A. J. *et al.* RADIA: RNA and DNA integrated analysis for somatic mutation
717 detection. *PLoS One* **9**, (2014).
- 718 51. Larson, D. E. *et al.* Somaticsniper: Identification of somatic point mutations in whole
719 genome sequencing data. *Bioinformatics* **28**, 311–317 (2012).
- 720 52. Koboldt, D. C. *et al.* VarScan: Variant detection in massively parallel sequencing of
721 individual and pooled samples. *Bioinformatics* **25**, 2283–2285 (2009).
- 722 53. Habegger, L. *et al.* VAT: a computational framework to functionally annotate variants in
723 personal genomes within a cloud-computing environment. *Bioinformatics* **28**, 2267–2269
724 (2012).
- 725 54. Smedley, D. *et al.* The BioMart community portal: an innovative alternative to large,
726 centralized data repositories. *Nucleic Acids Res.* **43**, W589-98 (2015).
- 727 55. Sethi, A., Eargle, J., Black, A. A. & Luthey-Schulten, Z. Dynamical networks in
728 tRNA:protein complexes. *Proc. Natl. Acad. Sci.* **106**, 6620–6625 (2009).
- 729 56. Girvan, M. & Newman, M. E. J. Community structure in social and biological networks.
730 *Proc. Natl. Acad. Sci.* **99**, 7821–7826 (2002).
- 731 57. Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R. & Siepel, A. Detection of nonneutral
732 substitution rates on mammalian phylogenies. *Genome Res.* **20**, 110–121 (2010).
- 733 58. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human
734 genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
- 735 59. The Gene Ontology Consortium. Expansion of the Gene Ontology knowledgebase and
736 resources. *Nucleic Acids Res.* **45**, D331–D338 (2017).
- 737 60. Fabregat, A. *et al.* The Reactome Pathway Knowledgebase. *Nucleic Acids Res.* **46**, D649–
738 D655 (2018).
- 739 61. Marisa, L. *et al.* KEGG: Kyoto encyclopedia of genes and genomes. *Nature* **10**, 1350–
740 1356 (2013).
- 741 62. Yu, G., Wang, L.-G., Han, Y. & He, Q.-Y. clusterProfiler: an R Package for Comparing
742 Biological Themes Among Gene Clusters. *Omi. A J. Integr. Biol.* **16**, 284–287 (2012).
- 743 63. Futreal, P. A. *et al.* A census of human cancer genes. *Nat. Rev. Cancer* **4**, 177–183 (2004).
- 744 64. Forbes, S. A. *et al.* COSMIC: Exploring the world’s knowledge of somatic mutations in
745 human cancer. *Nucleic Acids Res.* **43**, D805–D811 (2015).

- 746 65. Conway, J. R., Lex, A. & Gehlenborg, N. UpSetR: An R package for the visualization of
747 intersecting sets and their properties. *Bioinformatics* **33**, 2938–2940 (2017).
- 748 66. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new
749 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–
750 D361 (2017).
- 751 67. Forbes, S. A. *et al.* COSMIC: Somatic cancer genetics at high-resolution. *Nucleic Acids*
752 *Res.* **45**, D777–D783 (2017).
- 753 68. Miled, N. *et al.* The Phosphoinositide 3-Kinase Pathway. *Science (80-.)*. **296**, 1655–1657
754 (2007).
- 755 69. Davies, H. *et al.* Mutations of the BRAF gene in human cancer. *Nature* **417**, 949–954
756 (2002).
- 757 70. King, A. J. *et al.* Demonstration of a genetic therapeutic index for tumors expressing
758 oncogenic BRAF by the kinase inhibitor SB-590885. *Cancer Res.* **66**, 11100–5 (2006).
- 759 71. Veeriah, S. *et al.* The tyrosine phosphatase PTPRD is a tumor suppressor that is frequently
760 inactivated and mutated in glioblastoma and other human cancers. *Proc. Natl. Acad. Sci.*
761 *U. S. A.* **106**, 9435–40 (2009).
- 762 72. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-
763 associated genes. *Nature* **499**, 214–218 (2013).
- 764 73. Hoadley, K. A. *et al.* Cell-of-Origin Patterns Dominate the Molecular Classification of
765 10,000 Tumors from 33 Types of Cancer. *Cell* **173**, 291–304.e6 (2018).
- 766 74. Porta-Pardo, E. *et al.* Comparison of algorithms for the detection of cancer drivers at
767 subgene resolution. *Nat. Methods* **14**, 782–788 (2017).
- 768 75. Nogales, E. The development of cryo-EM into a mainstream structural biology technique.
769 *Nat. Methods* **13**, 24–7 (2016).
- 770
- 771
- 772
- 773