

1 **PhiDsc: Protein functional mutation Identification by 3D Structure Comparison**

2 Mohamad Hussein Hoballa¹ and Changiz Eslahchi^{1*}

3 ¹ Department of Computer Science, Shahid Beheshti University, Evin, Tehran, 1983963113 Iran

4 * To whom correspondence should be addressed. Tel: +98 21 22431653; Email: ch-eslahchi@sbu.ac.ir

5
6 Selective pressures that trigger cancer formation and progression shape the mutational landscape
7 of somatic mutations in cancer. Given the limits within which cells are regulated, a growing tumor
8 has access to only a finite number of pathways that it can alter. As a result, tumors arising from
9 different cells of origin often harbor identical genetic alterations. Recent expansive sequencing
10 efforts have identified recurrent hotspot mutated residues in individual genes. Here, we introduce
11 PhiDsc, a novel statistical method developed based on the hypothesis that, functional mutations in
12 a recurrently aberrant gene family can guide the identification of mutated residues in the family's
13 individual genes, with potential functional relevance. PhiDsc combines 3D structural alignment of
14 related proteins with recurrence data for their mutated residues, to calculate the probability of
15 randomness of the proposed mutation. The application of this approach to the RAS and RHO
16 protein families returned known mutational hotspots as well as previously unrecognized mutated
17 residues with potentially altering effect on protein stability and function. These mutations were
18 located in, or in proximity to, active domains and were indicated as protein-altering according to
19 six *in silico* predictors. PhiDsc is freely available at <https://github.com/hobzy987/PhiDSC-DALI>.

20

21

22 INTRODUCTION

23 Cancer development starts with the acquisition of genomic alterations and chromosomal
24 abnormalities that arise from uncorrected errors during DNA replication or repair or due to
25 exposure to mutagens (1). Some alterations may further the accumulation of somatic mutations (2)
26 and play a mechanistic role in malignant transformation. These “driver mutations” are postulated
27 to provide advantage to and promote cancer hallmarks in the subpopulation of cells that harbor
28 them (3). The number of driver mutations varies between cancer types, averaging four per tumor
29 (4). Most remaining somatic alterations, termed “passenger mutations,” may confer little to no
30 functional impact (5). However, distinguishing the handful of driver mutations from the vast
31 background of passenger mutations in a tumor has remained a challenge in cancer genomics.

32 Frequently altered nucleotides in the genes that are implicated in tumor development and
33 progression are known as mutational hotspots (6). The number of candidate hotspot mutations of
34 unknown functional significance has increased recently –especially due to the completion of large-
35 scale sequencing efforts such as The Cancer Genome Atlas (TCGA) (7), International Cancer
36 Genome Consortium (ICGC) (8), and Project GENIE (9). Many platforms are used to visualize
37 and organize these data like BioMuta (10) and cBioPortal(11, 12) allowing to download and
38 analyze large-scale cancer genomics datasets. Most of these frequently detected mutations are
39 within exons, or the coding regions of the proteins, and their function is ascertained by directly
40 examining their impact on the encoded protein or predicted through application of *in silico*
41 bioinformatic approaches (13, 14).

42 The statistical reoccurrence of mutations in tumors has been used as an indicator of their functional
43 impact, based on the assumption that infrequent alterations detected in tumors are likely non-
44 functional, passenger events (15). However, it has been shown that passenger mutations are not

45 randomly distributed along the cancer genomes (16). Rather, they are enriched in nucleotide
46 sequence contexts that are shaped by specific active mutational processes in a tumor (17, 18). In
47 contrast, driver mutations are postulated to occur in genomic positions whose distribution depends
48 not only on the local nucleotide context, but also on the location of functionally relevant residues
49 along the protein sequence (19, 20). Relying on recurrence alone to identify functional mutations,
50 may also be confounded by underlying mutational processes that target specific genomic contexts,
51 resulting in often-mutated residues that do not drive tumor progression (21).

52 In this context, numerous methods are presently being used to identify hotspot and driver
53 mutations, based on the frequency of mutations detected in a gene across a set of tumor samples
54 (e.g., MutSig (22) and MuSiC (23)). Recognizing mutational hotspot in infrequently altered genes
55 can also be refined by including protein-level annotation by local-positional clustering (24), or the
56 inclusion of phosphorylation sites (25) and information from paralogous protein domains (26).
57 Protein-level annotation, such as local-positional clustering, phosphorylation sites, and paralogous
58 protein domain (27) as well as 3D protein structures are used to identify functional mutations in
59 infrequently mutated genes.

60 Using a variety of approaches that take into account diverse aspects of protein structures and
61 types, functional mutations can be predicted across protein sequences and structures. Some
62 techniques, such as 3DHotspots (28), Hotspot3D (29), Mutation3D (30), and Signatures of
63 Cancer Mutation Hotspots in Protein Kinases (31) use the 3D structure of protein, while others
64 utilize 3D reconstruction of protein networks to provide a better understanding of genetic
65 abnormalities (32). On the other hand, methods like PinSnps (33), StructMAn (34), Hot-MAPS
66 (35) and SpacePAC (36), as well as SAAMBE-3D(37), use protein-protein interactions enriched
67 with somatic cancer mutations (38) to understand the effect of a mutation not only on the

68 function of the same protein but also on the signal transduction and activating cascade proteins.
69 Methods based on individual protein structures or the 3D reconstruction of protein networks
70 have improved the identification of mutational clusters in tumors ([39](#)) and have elucidated
71 functional consequences (folding free energy and stability of protein monomers ([40](#))) of protein-
72 altering mutations, other methods take into consideration the local DNA sequence context for the
73 analysis of cancer context-dependent mutations like MutaGene([41](#)). Although it is difficult to
74 categorize methods based on their input parameters (some require sequences while others may
75 need structures as well), in all cases, the output determines whether a proposed mutation has
76 occurred at a hotspot residue. However, a few limitations remain: First, focusing on the mutation
77 frequency across tumor samples increases the risk of missing portions of rare hotspot mutations
78 with low frequency; second, concentrating solely on driver genes fails to distinguish between
79 individual driver mutations within altered genes and passenger mutations within the same gene;
80 and third, analyzing protein sequences without a larger context misses the effect of mutations on
81 the conformational structure and functional sites of the protein.

82 To address these issues, we introduce PhiDsc. Its development is based on the hypothesis that
83 oncogenic mutations in a target protein can be identified by analyzing its three-dimensional
84 structural similarity, protein folding information, and mutational recurrence within its gene family.
85 We demonstrate that PhiDsc can identify candidate functional mutations, caused on altered protein
86 position, by comparing the three-dimensional structures of related human wild-type proteins and
87 assessing repeatedly altered residues in the protein family. PhiDsc combines the two approaches
88 by relying on the concept of hotspot mutations in functional regions and classifying protein
89 families based on their domains and active sites. Thus, by comparing the three-dimensional
90 structures of similar domains within a protein family, PhiDsc maps known functional mutations in
91 extensively studied proteins to those in the family that receive less interest.

92 **RESULTS**

93 PhiDsc is applied to HRAS from the RAS ([59](#)) subfamily and RhoA from the RHO ([60](#)) subfamily
94 of proteins.

95

96 **HRAS**

97 The family group of HRAS was $A(\text{HRAS}) = \{\text{DIRAS1}, \text{DIRAS2}, \text{GEM}, \text{KRAS}, \text{NRAS}, \text{RAP1A},$
98 $\text{RAP1B}, \text{RAP2A}, \text{RASL12}, \text{REM1}, \text{REM2}, \text{RERG}, \text{RRAD}, \text{RRAS}, \text{RRAS2}\}$. Dali aligned 98%
99 of HRAS residues to residues of each member of the family (**Table 1**) highlighting strong
100 structural similarity between the target protein and its respective protein families. (Supplementary
101 files HRAS alignment). As a result, PhiDsc scored 168 of 189 HRAS residues (89%) and predicted
102 13 residues as functional mutation (**Table 2**) all of which passed cross-validation evaluation

103 (Figure 1) and were consistently projected to be effective and protein-modifying by six
104 independent algorithms.

105 Table 1 indicates the percentage of structural alignment between each protein (HRAS) and its protein family member.

| Protein | HRAS | | | | | | | | | | | | |
|-----------|------|------|--------|--------|--------|--------|--------|--------|--------|--------|--------|------|--------|
| | RALA | RALB | RAP1A | RAP1B | RAP2A | KRAS | RASL12 | NRAS | RERG | RIT1 | RRAS2 | RRAS | Median |
| PDB ID | 2BOV | 2KWI | 1C1Y | 4DXA | 1KAO | 3GFT | 3CSC | 3CON | 2ATV | 4KLZ | 2ERY | 2FN4 | |
| Alignment | 100 | 100 | 97.619 | 98.214 | 98.214 | 98.809 | 95.238 | 92.857 | 98.809 | 92.261 | 97.619 | 100 | 98.214 |

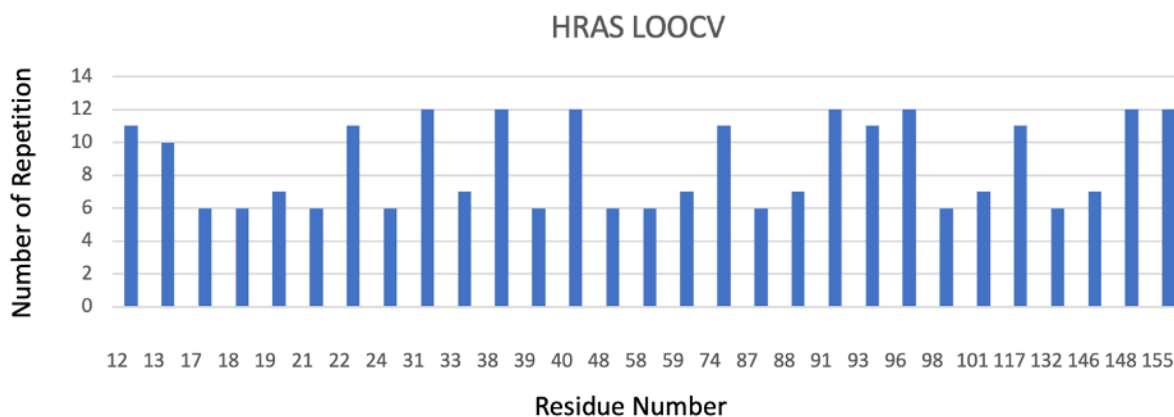
106

107 Table 2 Candidate functional mutations for HRAS proposed by PhiDsc. Residue positions sorted by their PhiDsc score p-value along
108 with predicted interacting residues from the RIN analysis are shown. COSMIC mutation reference or dbSNP polymorphism ID are

109 also shown when available.

| HRAS | | | | | | | | | | | |
|------------|----------|---------------------|-----|-----|-----|-----|-----|-----|-----|-----|-------------|
| Residue Nu | P-value | Interacting Residue | | | | | | | | | mut ref Nu |
| 12 | 2.66E-07 | 11 | 16 | | | | | | | | COSM483 |
| 74 | 1.72E-06 | 5 | 70 | 71 | 73 | 75 | | | | | COSM5991570 |
| 13 | 2.57E-06 | 117 | | | | | | | | | COSM486 |
| 93 | 6.18E-06 | 81 | 82 | 90 | 91 | 113 | 137 | | | | COSM9497546 |
| 91 | 8.72E-06 | 87 | 88 | 90 | 93 | 95 | | | | | COSM6476473 |
| 22 | 1.38E-05 | 18 | 19 | 20 | 32 | 26 | 28 | 146 | 149 | 152 | COSM6923245 |
| 96 | 1.54E-05 | 9 | 10 | 11 | 92 | 93 | 97 | 98 | 99 | 100 | RS889495169 |
| 117 | 1.85E-05 | 13 | 14 | 83 | 84 | 116 | 119 | 120 | 144 | | CSOM304967 |
| 31 | 3.84E-05 | 30 | 33 | | | | | | | | COSM6915342 |
| 40 | 4.20E-05 | 20 | 24 | 32 | 38 | 39 | 54 | 55 | 57 | | RS763920334 |
| 155 | 5.08E-05 | 79 | 144 | 151 | 152 | 153 | 159 | | | | COSM9515051 |
| 148 | 5.23E-05 | 119 | 145 | 150 | | | | | | | COSM6903495 |
| 38 | 5.93E-05 | 39 | 40 | 57 | | | | | | | RS750680771 |

110



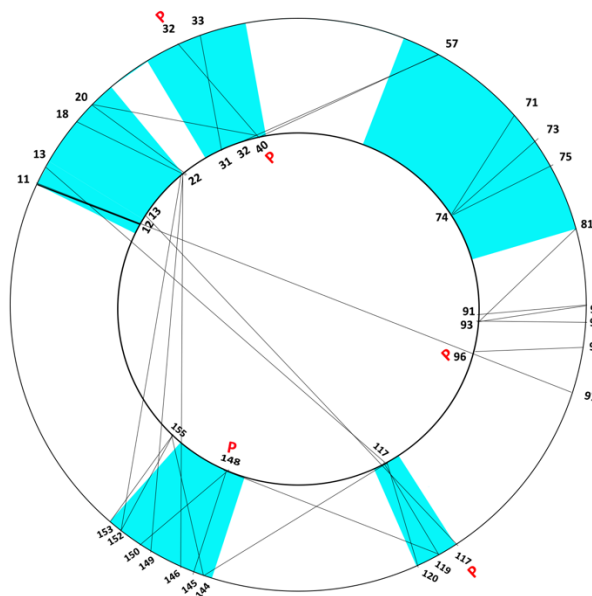
111

112 Figure 1 shows the LOOCV for the two proteins HRAS, in all the iterations of the system the number of repeated times for each
 113 residue is shown, (>80%), which indicates that the results obtained by the system are robust since the original results are obtained
 114 in all the LOOCV iteration

115 RIN is generated using the HRAS structure (RCSB database ID 4Q21, with 168 residues). Thirteen
 116 candidate functional mutations shared 58 neighboring residues located in the functional domains
 117 of the protein (G boxes, Switches I and II, GDI and GEF interaction sites, GTP/MG2+ binding
 118 domain). Moreover, 25 of these 58 residues were seen mutated in human tumors according to the
 119 cBioPortal([11](#), [12](#)) database a distinct dataset form BioMuta.

120 Top-four PhiDsc predictions in HRAS were residues 12, 13, 74, and 93, which are known to be
121 key functionals and often mutated in various cancer types (61). The domain comprising residues
122 12 and 13 is involved in Guanine Nucleotide Dissociation Inhibitor (GDI) interaction as well as
123 interaction with GTP/Mg²⁺ (62), and is mostly detected in tumors such as bladder cancer (63),
124 thyroid cancer(64), and other diseases such as Costello syndrome (61) and Schimmelpenning-
125 Feuerstein-Mims syndrome (63, 65). Mutations in residue 74 are seen in endometrioid cancer and
126 sebaceous carcinoma, while those in residue 93, have been discovered in only a small percentage
127 of prostate cancer samples (66). According to Ensemble Learning Approach for Stability
128 Prediction of Interface and Core mutations (ELSPIC) (67), residue 93 is localized in the protein's
129 core, suggesting that it has a direct effect on the protein's shape and function.

130 Although 3 of 13 candidate functional mutations in HRAS were not located in any protein domains,
131 they were found near the intersection of exons 3 and 4 at residue 97. Finally, residue 96 has been
132 identified as a phosphorylation site, the other residues as shown in (Figure 2) were located in
133 functional protein domains.



134

135 *Figure 2 depicts the inner circle's candidate functional mutations and the outer circle's interacting residues. According to*
136 *thecanSAR BLACK system (60), The blue areas represent HRAS functional regions, while the lines linking the inner circle (candidate*
137 *functional mutation) to the outer circle (interacting residues) represent residue interactions. This figure displays only the*
138 *HRAS residues that are mutated in cBioPortal.*

139 RhoA

140 RhoA, a member of the RHO (60) subfamily of proteins with $A(\text{RhoA}) = \{\text{RHOB}, \text{RHOC}, \text{RHOD},$
141 $\text{RHOQ}, \text{RHOU}, \text{RND1}, \text{RND3}, \text{RAC1}, \text{RAC2}, \text{RAC3}, \text{CDC42}\}.$

142 The RCSB database is used to retrieve 3D structure files for each member (if found in PDB)
143 of A(RhoA). The final list of PDB structures are shown in **Table 3**. The Dali server is then used
144 to perform a pairwise structural comparison between the input protein and each member of its
145 family. 97% of RhoA residues were aligned with the residues of each family member in the
146 generated alignments. The existence of strong structural similarities between target proteins and
147 their respective protein families supports these results (Supplementary file “RhoA alignment”).

148 As an outcome, 179 out of 193 residues were scored for RhoA.

149 *Table 3 shows the percentage of structural alignment of each protein (RhoA) with its corresponding protein family member.*

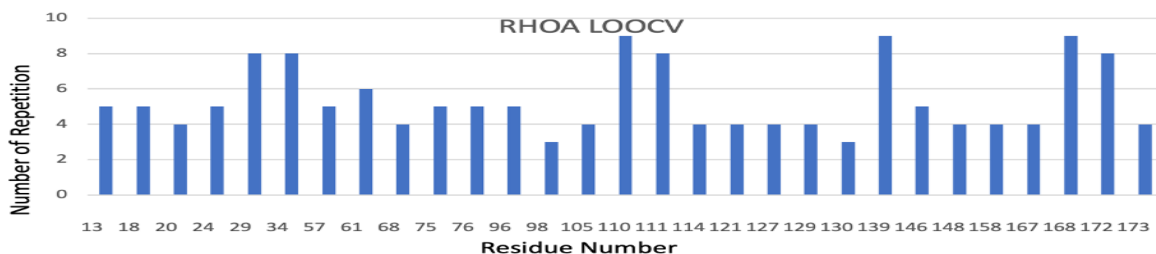
| RHOA | | | | | | | | | | | |
|-----------|--------|--------|--------|-------|--------|--------|--------|--------|--------|------|--------|
| Protein | RAC1 | RAC2 | RAC3 | RHOB | RHOC | RHOD | RHOQ | RHOV | RND1 | RND3 | Median |
| PDB ID | 1E96 | 1DS6 | 2C2H | 2FV8 | 2GCN | 2J1L | 2ATX | 2Q3H | 3Q3J | 2V55 | |
| Alignment | 99.441 | 99.441 | 93.854 | 95.53 | 98.324 | 87.709 | 99.441 | 94.413 | 97.206 | 100 | 97.765 |

150
 151 The P-value of PhiDsc statistics is generated for all target protein residues in the final phase. Eight
 152 candidate functional mutations for RhoA were obtained. **Table 4** illustrates the RhoA protein
 153 candidate functional mutations introduced by the PhiDsc procedure. The eight candidates passed
 154 cross validation (see **Figure 3**) and were consistently predicted to be effective and protein-
 155 modifying by six separate algorithms. Despite the fact that no evidence of a mutation in residue
 156 29 of RhoA was detected in any cancer mutation databases, all six techniques predicted that this
 157 mutation would alter RhoA's functional activity.

158 *Table 4 lists all candidate functional mutations for RhoA proposed by the PhiDsc approach. The table shows the residue position*
 159 *number (P) in the first column, sorted by their P-value in the second column, the interacting residues of each candidate functional*
 160 *mutation in the third column, the "COSM" letters of the mutations indicate that these mutations were annotated in the cosmic*
 161 *database as tumor-related mutations, while the "rs" letters of the mutations indicate that these mutations were annotated in the*
 162 *Dpsnp database.*

| RHOA | | | | | | | | | | | | |
|----------------|-------------|---------------------|-----|-----|-----|-----|-----|-----|-----|--|--|-----------------|
| Residue Number | P-value | interacting residue | | | | | | | | | | mutation ref NU |
| 111 | 3.07904E-05 | 78 | 79 | 80 | 109 | 110 | 177 | | | | | COSM2849881 |
| 34 | 4.86819E-05 | 35 | | | | | | | | | | COSM2849895 |
| 139 | 9.956E-05 | 84 | 86 | 89 | 92 | 122 | 139 | 140 | 143 | | | COSM2849897 |
| 168 | 0.000147858 | 170 | 171 | 172 | | | | | | | | COSM7114068 |
| 110 | 0.000209526 | 77 | 78 | 79 | 80 | 107 | 108 | 11 | | | | RS368767616 |
| 29 | 0.000224094 | 23 | 27 | 28 | 29 | 31 | | | | | | NO |
| 172 | 0.000300752 | 46 | 48 | 168 | 169 | 172 | 174 | 175 | 176 | | | COSM1309264 |
| 127 | 0.000484266 | 87 | 121 | 124 | 125 | 127 | 129 | 130 | 131 | | | MU85445108 |

163



164

165 *Figure 3 shows the LOOCV for the protein RhoA; the number of repeated times for each residue is presented in all iterations of the*
166 *system, indicating that the system's results are resilient because the original results are obtained in all LOOCV iterations.*

167 The RIN for RhoA is constructed using 1OW3 obtained from the RCSB database. The 8

168 potential functional mutations have 42 neighbors, 18 of which had previously been identified as

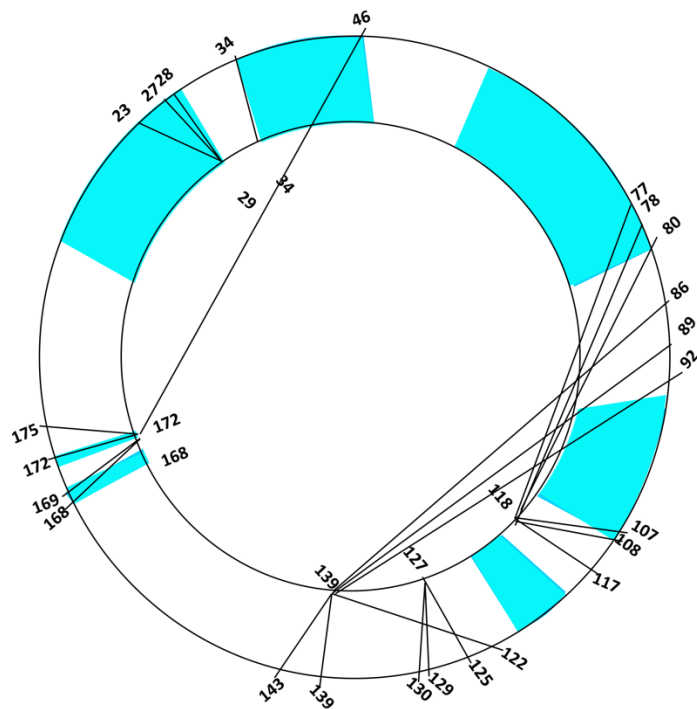
169 occurring mutations in the cBioPortal database ([11](#), [12](#)) (see **Table 3/interacting residues**). The

170 neighbors of potential functional mutations are related to PPI functionals, according to

171 RINalyzer data. These neighbors are also located in RhoA protein domains associated to GAP,

172 GEF, and GDI interaction and phosphorylation sites, including position 127—showing that this

173 residue is significant in RhoA's functional activity (see **Figure 4**).



174

175 *Figure 4 shows the inner circle's candidate functional mutations and the outer circle's interacting residues. According to*
176 *thecanSAR BLACK system (60), The blue areas represent RhoA functional regions, while the lines linking the inner circle*
177 *(candidate functional mutation) to the outer circle (interacting residues) represent residue interactions. This figure displays only*
178 *the HRAS residues that are mutated in cBioPortal.*

179 In cancer samples, four high-scoring RhoA residues (34, 139, 111, and 168) were observed (see
180 **Table 4**). Residue 34 is near the core area and the GAP interaction site, as per RhoA's 3D structure.
181 A mutation at this location improves the affinity for ARHGAP;1, a GAP protein that plays a vital
182 role in RhoA activation, according to data from ELASPIC (67) and COSMIC. According to
183 COSMIC, mutation 139 of RhoA was observed in one sample of non-small cell lung carcinoma
184 and as a silent mutation in two samples of cervix and stomach cancer— where it was not a
185 functional mutation in the latter two samples. Meanwhile, residue 111 has been seen in one sample
186 of stomach cancer patients (7). Mutation in residue 168 boosts the affinity for the CTRO protein,
187 which regulates cytokinesis by generating a contractile ring. It was also found to interact with
188 KAPCA, a gene associated with breast and ovarian cancer (68). The mutation of residue 168 also

189 impacted PKN1 and PKN2 interaction with RhoA—two proteins that contribute to prostate cancer
190 and play a crucial role in cell migration and proliferation ([69](#), [70](#)).

191 **DISCUSSION**

192 In this paper, we looked at proteins that are similar and have been classified into families in
193 uniprotkb. In terms of sequence, structure, and function, these proteins are very similar. As a result,
194 we assume that the frequent mutations associated with the same cancer phenotype on the same
195 domain share these domains and mutations within the family. As a result, the introduced algorithm
196 employs scores to determine whether these mutations are statistically significant as functional
197 alterations in areas common in families. To test and validate the approach, domains from two well-
198 known protein families (HRAS and RhoA) that are known to be involved in cancer are used.

199 As a result, we present PhiDsc, a novel method for detecting functional mutations in proteins. To
200 link mutation residues to specific biological functional domains of proteins, we took into account
201 a mutation's position in the protein's 3D structure ([71](#)), as well as the frequency of its reoccurrence
202 in human tumors ([72](#)). Finally, we combined these characteristics with known functional hotspot
203 mutations aggregated among paralogous proteins in the same family or with similar domains ([73](#)),
204 and we used Bonferroni restriction to further narrow the range of predictions in order to reduce
205 false positives..

206

207 We evaluated PhiDsc using the HRAS and RhoA proteins ([71](#), [72](#)). HRAS is a GTPase protein in
208 the RAS subfamily that controls many cellular mechanisms including 84 pathways according to
209 KEGG Pathway. The most mutated residues in HRAS are 12, 13, and 61, which are related to
210 different subsets in cancer ([73](#)), and the tumorigenic effect of HRAS is related to the protein's

211 permanent activation. RhoA is a RHO subfamily signaling G protein that regulates numerous
212 cellular mechanisms associated with 43 pathways related to cellular processes as seen in KEGG
213 Pathway. The most frequently mutated residues in this protein, 17 and 42, have been observed in
214 various types of cancer (74), and similarly, the oncogenic effect of RhoA is exerted by its constant
215 activation of the protein.

216 With the exception of one candidate residue in RhoA, all residues predicted by PhiDsc were found
217 to be mutated in cancer samples, as well as in other diseases such as Costello syndrome, which is
218 linked to germline HRAS mutations (75). Although certain candidate functional mutations were
219 not previously identified as hotspot mutations and had a low mutated frequency in cancer mutation
220 datasets (rare), using CanSar balck (60), we demonstrated that they were located in active
221 functional domains of proteins or had a wide network of interactions with functional residues.
222 Noteworthy, the Biomuta database was initially used; however, by the final step, some of the
223 candidate functional hotspots that were not found in Biomuta had been presented in tumor samples
224 in other datasets such as COSMIC, cBioPortal, and Dbsnp. With the exception of RhoA residue
225 29, all were identified as rare mutated residues, and, thus, they were not previously mentioned as
226 a hotspots, indicating that PhiDsc improves and optimizes the detection of low frequency
227 functional mutations. while, residue 29 of RhoA had no mutational records in COSMIC (46) or
228 Dbsnp (76) databases, mutation analysis software MutaGene (41) ranked RhoA residue 29 as a
229 highly mutable position, and the projected effect by six different software packages at that position
230 predicts a potential oncogenic effect. It is notable that the difference between COSMIC and Dbsnp
231 lies in the curation method used to classify any given mutation as an SNP.

232 Despite the fact that these methods use different concepts to infer the stabilizing effect of point
233 mutations (as discussed in the results section), they all suggest that PhiDsc's predictions alter

234 protein structure and function. The precise impact of unknown mutations necessitates additional
235 experimental verification.

236 When DALI was used instead of TM-Align, better results were obtained in PhiDsc with known
237 functional mutations. These findings suggest that different 3D alignment approaches may alter
238 predicting hotspot mutations in different types of proteins. As a result, the PhiDsc package's
239 predictions should improve as the mode of alignment used improves.

240 Some previously designated hotspots of HRAS and RhoA in cancer, like for HRAS out of 12
241 (residues 12, 13 and 117) and for RhoA out of 11 (residue 34) were returned by PhiDsc. When the
242 results of the Dali and Tm-Align alignment (supplementary files (HRAS, RHOA) Tm-Align)
243 methods were compared, the results of the Tm-Alignment method predicted fewer well-known
244 driver mutations than the results of the Dali method. This suggests that a different alignment choice
245 could result in some differences in predictions.

246 Although the two example proteins selected for validation are oncogenic, PhiDsc is not restricted
247 to oncogenes and can be utilized to identify functional mutations in tumor suppressor genes or any
248 other type of Protein if the family has a sufficient number of members and the mutation profile
249 data is adequate and consistent.

250 The lack of a 3D structure of the protein and small protein families, which limit the number of
251 members in the family, are two limitations of this method. A future update to the tool will
252 include the ability to align functional domains of proteins rather than the entire protein, as well
253 as the use of the protein's predicted 3D structure in the alignment comparison.

254 MATERIALS AND METHODS

255 PhiDsc Algorithm

256 PhiDsc uses a six-step method that is centered on a protein **P** with **m** amino acid residues and a
257 known three-dimensional structure. Briefly, a list of proteins is defined, denoted by the set **A(P)**,
258 by identifying all members of P's protein family from UniProtKB (42) and selecting all human
259 proteins with 3D structures from the Protein Data Bank (PDB) (43). Next, the 3D structures of the
260 proteins members in A(P) are aligned to the 3D structure of P. The results are presented by a
261 matrix, **E(P)**. Then, using the BIOMUTA V4 and 3Dhotspot database (44), the mutational
262 information of each protein of A(P) is identified, in order to score each residue of P and calculate
263 an associated probability. Finally, these are analyzed to identify potential candidate functional
264 mutations in P. Each step is described in detail in what follows.

265 Step 1: Define the protein list A(P). The UniProtKB database (42) is used to identify members of
266 a given protein's protein family, while the RCSB Protein Data Bank (PDB) (43) is used to
267 determine their three-dimensional structure. The PDB contains the structures of wild-type and
268 mutated proteins. For the alignment step, either the full-length sequence of the wild-type protein
269 or the least mutated form (maximum one mutation) of the same length is used; the final list is
270 denoted by $A(P) = \{P_1, P_2, P_3 \dots P_n\}$.

271 Step 2: Align 3D structures. Dali, a pairwise comparison server for protein structures, is used to
272 align protein structures (<http://ekhidna2.biocenter.helsinki.fi/dali/>)(45). TM-Align “another
273 alignment method” is also included in PhiDsc with its default parameters.

274 Step 3: Define matrix E(P). $E(P) = [a_{kij}^j]$ has n columns (number of proteins) and m rows (number
275 of amino acids in protein P), in which a_{kij}^j denotes the type of amino acid in the sequence of protein

276 j that is aligned to the i^{th} amino acid in protein P; k_{ij} denotes the position number of amino acid in
277 the sequence P_j that is aligned to the i^{th} amino acid of protein P.

278 Step 4: Identify mutational information of each protein in A(P). Residues for all protein family
279 members are annotated with mutational and hotspot information using BioMuta (version 4, (10))
280 and 3Dhotspots (39). BioMuta is a database of curated cancer-associated single-nucleotide
281 variations derived from COSMIC (46), ClinVar (47), CIVIC(48), and UniProtKB(42) and actively
282 curated from publications and automated analysis of publicly available databases such as
283 TCGA(7)and ICGC(8). 3Dhotspots is a dataset of statistically significant mutations clustered in
284 three-dimensional protein structures found in cancer. The data set contains mutational positions
285 referred to as hotspot mutations.

286 Step 5: Score residues. A grade is assigned to each amino acid of A(P) members based on the
287 mutational information for that amino acid (P). Let a_k^t be the k th amino acids of protein P_t . Define:

$$288 \quad m(a_k^t) = \begin{cases} 1, & \text{if } a_k^t \text{ is reported as mutation in biomuta} \\ 2, & \text{if } a_k^t \text{ is reported as hotspot in 3Dhotspots database} \\ 0, & \text{otherwise (either non - aligned or not mutated)} \end{cases}$$

289

290 Let the i^{th} row of the matrix E(P) be $[a_{k_{i1}}^1, a_{k_{i2}}^2, \dots, a_{k_{in}}^n]$, $1 \leq i \leq m$. The following score is
291 assigned to i^{th} amino acids of P:

$$292 \quad S(i) = \sum_{j=1}^n m(a_{k_{ij}}^j)$$

293 To calculate the statistical significance of the obtained scores $S(i)$ at each position (row in the
294 matrix $E(P)$), we calculate the probability related to this score. Let protein P_t have m_t amino acids
295 of which l_t are mutated in biomuta. Define:

$$296 \quad P(a_k^t) = \begin{cases} \frac{l_t}{m_t}, & m(a_k^t) > 0 \\ 1 - \frac{l_t}{m_t}, & m(a_k^t) = 0 \end{cases}$$

297 To distinguish non-mutated from the non-aligned residues (both with score $m(a_k^t) = 0$), and
298 because the event under investigation is the occurrence of functional mutations that are coded in
299 the alignments. Then, if in $a_{k_{ij}}^j$ (j) is a gap, we assume $P(a_{k_{ij}}^j) = 1$.

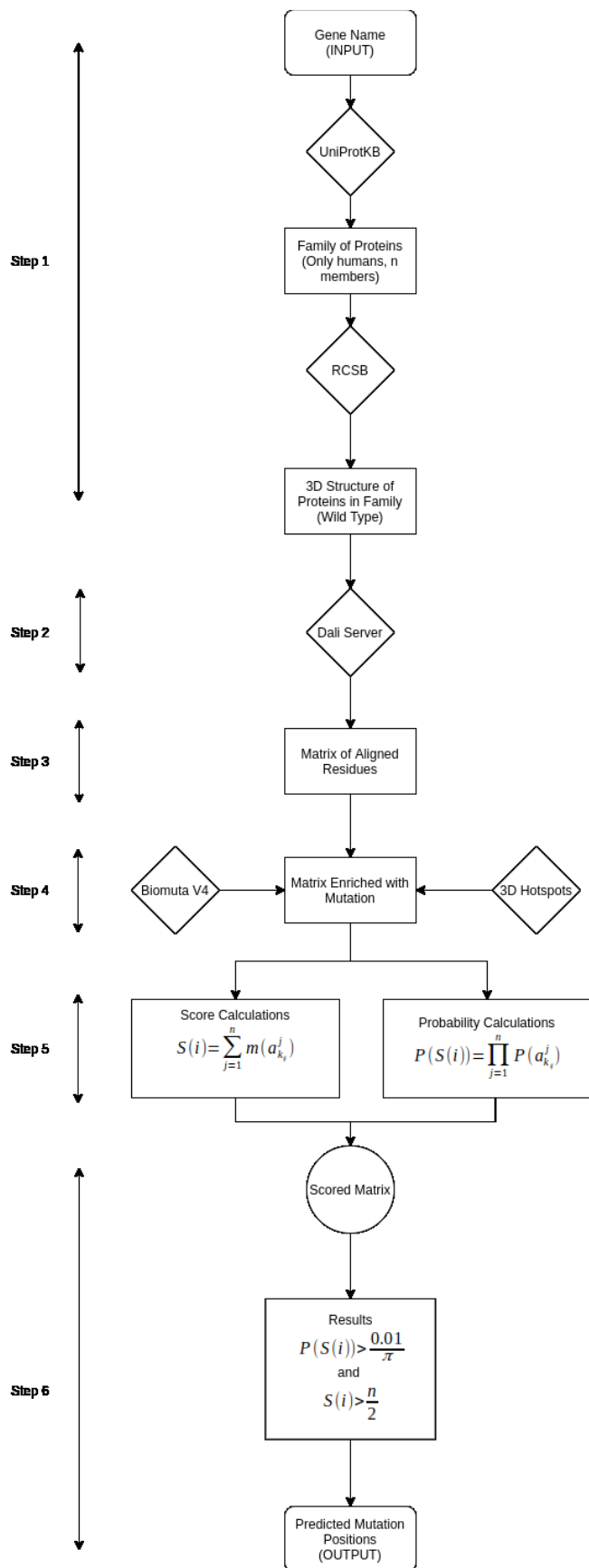
300 Then:

$$301 \quad P(S(i)) = \prod_{j=1}^n P(a_{k_{ij}}^j)$$

302

303 Step 6: Select candidates. The i^{th} amino acid of protein P is selected as a candidate functional
304 mutation if $P(S(i))$ is less than $\frac{0.01}{n}$, following the Bonferroni correction, and if $S(i) > \frac{n}{2}$.

305 The method is schematically described in **Figure 5**



307 *Figure 5 The PhiDsc workflow. The system begins by obtaining family members; the algorithm then obtains the 3D structures from*
308 *RCSB; the algorithm aligns members pairwise with the input protein; mutations are then enriched in the alignments; finally, scores*
309 *and probabilities are calculated.*

310

311 **Leave-one-out cross-validation**

312 In leave-one-out cross-validation (LOOCV), one data point from the training set remains excluded.

313 For example, if there are n data points in the original sample, $n-1$ samples are used to train the

314 model, and p points are used as the validation set. This is repeated for all combinations in which

315 the original sample can be separated in this manner, and the error is averaged across all trials to

316 calculate overall effectiveness. The number of possible combinations is equal to the original

317 sample's number of data points, or n .

318 $A_i(P) = \{P_1, P_2, P_3 \dots P_n\} - \{P_i\}$ is considered as an input set for protein P , and the PhiDsc

319 predictions for P are obtained by considering $A_i(P)$ as its protein family set. The set of predicted

320 functional mutations is obtained for every $1 \leq i \leq n$. A projected functional mutation is said to

321 be robust if it is predicted across at least 80% of all rounds.

322 **Residue Interaction Network**

323 RIN (Residue Interaction Network) is used to quantify the physical effect of the mutation on

324 protein structure and function. In summary, Chang *et al.* demonstrated that if a mutation in a

325 protein's 3D structure is close to some hotspot mutations, the likelihood of this mutation being

326 considered a hotspot mutation is high. The RINalyzer ([49](#)) module generates user-defined RINs

327 from a 3D protein structure obtained from RCSB protein databank. RINerator considers different

328 biochemical interaction types, such as contacts/clashes, hydrogen bonds, and hydrogen atoms and

329 quantifies their individual strength as described in Chimera ([50](#)). RINalyzer is a Java plugin for

330 Cytoscape([51](#)), a free software platform for the analysis and visualization of molecular interaction

331 networks. The results of interacting residues from RIN are compared to cBioPortal ([11](#), [12](#)) a
332 dataset of mutations that are curated across cancer samples.

333 **Functional effect of candidate mutations on proteins**

334 The effect of alterations in regions that were not identified as functional mutations experimentally
335 can be calculated using a variety of methods. PhiDsc's functional predictions are evaluated using
336 six methods that, according to Stefl et al. ([52](#)), can be classified into three types:

337 The first group includes machine learning approaches that are trained on protein stability features
338 and account for experimental conditions such as temperature, salt concentration, and pH values.
339 Incorporating such parameters is critical for assessing the free-energy changes caused by mutations
340 under near physiological conditions. This group includes I-Mutant2.0 ([53](#)) which uses SVM to
341 estimate $\Delta\Delta G$ upon mutation, and PoPMuSiC-2.0 ([54](#)) which uses a mix of statistical potential and
342 neural networks to estimate $\Delta\Delta G$ upon mutation.

343 The second group relies on evolutionary conservation data, with the assumption that changes at
344 conserved positions in multiple sequence alignments are detrimental. Although these approaches
345 do not directly predict the effect of mutations on protein stability, they are commonly used in
346 conjunction with the methods mentioned above to achieve consensus predictions. This group
347 includes SIFT ([55](#)), which uses sequence homology and site conservation to estimate the
348 deleterious effect of mutations, and Provean ([56](#)), which predicts the functional impact of all types
349 of protein sequence variations, including single amino acid substitutions, insertions, deletions, and
350 multiple substitutions.

351 The third group uses structural information, assuming that a protein's ability to function properly
352 is determined by fundamental physicochemical properties that can only be derived from structures.

353 This group includes CUPSAT([57](#)), which estimates $\Delta\Delta G$ upon mutation using mean force atom
354 pair and torsion angle potentials, and MutPred([58](#)), which estimates detrimental effect of mutation
355 using SIFT and gain/loss of structural or functional features predicted from sequences.

356

357 **DATA AVAILABILITY**

358 This method is implemented in Python and the Source code and all tested data can be found on
359 (<https://github.com/hobzy987/PhiDSC-DALI>). The software takes a UniProt Protein name as
360 input and gives html file as output with aligned residues and probabilities, and a list of all residues
361 sorted according to their score.

362 **ACKNOWLEDGEMENT**

363 The authors thank Dr. Hossein Khiabani for the insightful discussions and contribution provided
364 for this work.

365

366 **REFERENCES**

- 367 1. Vogelstein B, Papadopoulos N, Velculescu VE, Zhou S, Diaz LA, Jr., Kinzler KW. Cancer
368 genome landscapes. *Science*. 2013;339(6127):1546-58.
- 369 2. Nowell PC. The clonal evolution of tumor cell populations. *Science*. 1976;194(4260):23-8.
- 370 3. Hanahan D, Weinberg RA. Hallmarks of cancer: the next generation. *cell*. 2011;144(5):646-74.
- 371 4. Martincorena I, Raine KM, Gerstung M, Dawson KJ, Haase K, Van Loo P, et al. Universal
372 patterns of selection in cancer and somatic tissues. *Cell*. 2017;171(5):1029-41. e21.
- 373 5. Stratton MR, Campbell PJ, Futreal PA. The cancer genome. *Nature*. 2009;458(7239):719-24.

- 374 6. Baeissa H, Benstead-Hume G, Richardson CJ, Pearl FM. Identification and analysis of mutational
375 hotspots in oncogenes and tumour suppressors. *Oncotarget*. 2017;8(13):21290.
- 376 7. Tomczak K, Czerwińska P, Wiznerowicz M. The Cancer Genome Atlas (TCGA): an
377 immeasurable source of knowledge. *Contemporary oncology*. 2015;19(1A):A68.
- 378 8. Zhang J, Bajari R, Andric D, Gerthoffert F, Lepsa A, Nahal-Bose H, et al. The international
379 cancer genome consortium data portal. *Nature biotechnology*. 2019;37(4):367-9.
- 380 9. Consortium APG. AACR Project GENIE: powering precision medicine through an international
381 consortium. *Cancer discovery*. 2017;7(8):818-31.
- 382 10. Dingerdissen HM, Torcivia-Rodriguez J, Hu Y, Chang T-C, Mazumder R, Kahsay R. BioMuta
383 and BioXpress: mutation and expression knowledgebases for cancer biomarker discovery. *Nucleic acids
384 research*. 2018;46(D1):D1128-D36.
- 385 11. Cerami E, Gao J, Dogrusoz U, Gross BE, Sumer SO, Aksoy BA, et al. The cBio cancer genomics
386 portal: an open platform for exploring multidimensional cancer genomics data. *AACR*; 2012.
- 387 12. Gao J, Aksoy BA, Dogrusoz U, Dresdner G, Gross B, Sumer SO, et al. Integrative analysis of
388 complex cancer genomics and clinical profiles using the cBioPortal. *Science signaling*. 2013;6(269):p11-
389 pl.
- 390 13. Martelotto LG, Ng CK, De Filippo MR, Zhang Y, Piscuoglio S, Lim RS, et al. Benchmarking
391 mutation effect prediction algorithms using functionally validated cancer-related missense mutations.
392 *2014;15(10):1-20*.
- 393 14. Tchernitchko D, Goossens M, Wajcman HJCC. In silico prediction of the deleterious effect of a
394 mutation: proceed with caution in clinical genetics. *2004;50(11):1974-8*.
- 395 15. Taylor BS, Barretina J, Socci ND, DeCarolis P, Ladanyi M, Meyerson M, et al. Functional copy-
396 number alterations in cancer. *PloS one*. 2008;3(9):e3179.
- 397 16. Dietlein F, Weghorn D, Taylor-Weiner A, Richters A, Reardon B, Liu D, et al. Discovery of
398 cancer driver genes based on nucleotide context. *bioRxiv*. 2018:485292.

- 399 17. Alexandrov LB, Nik-Zainal S, Wedge DC, Aparicio SA, Behjati S, Biankin AV, et al. Signatures
400 of mutational processes in human cancer. *Nature*. 2013;500(7463):415-21.
- 401 18. Nik-Zainal S, Davies H, Staaf J, Ramakrishna M, Glodzik D, Zou X, et al. Landscape of somatic
402 mutations in 560 breast cancer whole-genome sequences. *Nature*. 2016;534(7605):47-54.
- 403 19. Chakravorty D, Jana T, Mandal SD, Seth A, Bhattacharya A, Saha S. MYCbase: a database of
404 functional sites and biochemical properties of Myc in both normal and cancer cells. *BMC bioinformatics*.
405 2017;18(1):1-10.
- 406 20. Chang MT, Bhattarai TS, Schram AM, Bielski CM, Donoghue MT, Jonsson P, et al. Accelerating
407 discovery of functional mutant alleles in cancer. *Cancer discovery*. 2018;8(2):174-83.
- 408 21. Makova KD, Hardison RCJNRG. The effects of chromatin organization on variation in mutation
409 rates in the genome. 2015;16(4):213-23.
- 410 22. Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K, Sivachenko A, et al. Mutational
411 heterogeneity in cancer and the search for new cancer-associated genes. *Nature*. 2013;499(7457):214-8.
- 412 23. Phan DL, Kim Y, Kim M, editors. MUSIC: Mutation analysis tool with high configurability and
413 extensibility. 2018 IEEE International Conference on Software Testing, Verification and Validation
414 Workshops (ICSTW); 2018: IEEE.
- 415 24. Tamborero D, Gonzalez-Perez A, Lopez-Bigas NJB. OncodriveCLUST: exploiting the positional
416 clustering of somatic mutations to identify cancer genes. 2013;29(18):2238-44.
- 417 25. Reimand J, Bader GDJMSb. Systematic analysis of somatic mutations in phosphorylation
418 signaling predicts novel cancer drivers. 2013;9(1):637.
- 419 26. Miller ML, Reznik E, Gauthier NP, Aksoy BA, Korkut A, Gao J, et al. Pan-cancer analysis of
420 mutation hotspots in protein domains. 2015;1(3):197-209.
- 421 27. Cantor AJ, Shah NH, Kuriyan J. Deep mutational analysis reveals functional trade-offs in the
422 sequences of EGFR autophosphorylation sites. *Proceedings of the National Academy of Sciences*.
423 2018;115(31):E7303-E12.

- 424 28. Gao J, Chang MT, Johnsen HC, Gao SP, Sylvester BE, Sumer SO, et al. 3D clusters of somatic
425 mutations in cancer reveal numerous rare mutations as functional targets. 2017;9(1):1-13.
- 426 29. Chen S, He X, Li R, Duan X, Niu B. HotSpot3D web server: an integrated resource for mutation
427 analysis in protein 3D structures. *Bioinformatics*. 2020;36(12):3944-6.
- 428 30. Meyer MJ, Lapcevic R, Romero AE, Yoon M, Das J, Beltrán JF, et al. mutation3D: cancer gene
429 prediction through atomic clustering of coding variants in the structural proteome. *Human mutation*.
430 2016;37(5):447-56.
- 431 31. Chen W, Li Y, Wang Z. Evolution of oncogenic signatures of mutation hotspots in tyrosine
432 kinases supports the atavistic hypothesis of cancer. *Scientific reports*. 2018;8(1):1-8.
- 433 32. Wang X, Wei X, Thijssen B, Das J, Lipkin SM, Yu H. Three-dimensional reconstruction of
434 protein networks provides insight into human genetic disease. *Nature biotechnology*. 2012;30(2):159-64.
- 435 33. Lu H-C, Herrera Braga J, Fraternali FJB. PinSnps: structural and functional analysis of SNPs in
436 the context of protein interaction networks. 2016;32(16):2534-6.
- 437 34. Gress A, Ramensky V, Büch J, Keller A, Kalinina OV. StructMAN: annotation of single-
438 nucleotide polymorphisms in the structural context. *Nucleic acids research*. 2016;44(W1):W463-W8.
- 439 35. Tokheim C, Bhattacharya R, Niknafs N, Gygyax DM, Kim R, Ryan M, et al. Exome-Scale
440 Discovery of Hotspot Mutation Regions in Human Cancer Using 3D Protein Structure.
441 2016;76(13):3719-31.
- 442 36. Ryslik G, Cheng Y, Zhao H. SpacePAC: Identifying mutational clusters in 3D protein space
443 using simulation. 2013.
- 444 37. Pahari S, Li G, Murthy AK, Liang S, Fragoza R, Yu H, et al. SAAMBE-3D: Predicting effect of
445 mutations on protein-protein interactions. 2020;21(7):2563.
- 446 38. Wong ET, So V, Guron M, Kuechler ER, Malhis N, Bui JM, et al. Protein-protein interactions
447 mediated by intrinsically disordered protein regions are enriched in missense mutations. *Biomolecules*.
448 2020;10(8):1097.

- 449 39. Gao J, Chang MT, Johnsen HC, Gao SP, Sylvester BE, Sumer SO, et al. 3D clusters of somatic
450 mutations in cancer reveal numerous rare mutations as functional targets. *Genome medicine*. 2017;9(1):1-
451 13.
- 452 40. Banerjee A, Mitra P. Estimating the Effect of Single-Point Mutations on Protein Thermodynamic
453 Stability and Analyzing the Mutation Landscape of the p53 Protein. *Journal of chemical information and*
454 *modeling*. 2020;60(6):3315-23.
- 455 41. Goncarenco A, Rager SL, Li M, Sang QX, Rogozin IB, Panchenko AR. Exploring background
456 mutational processes to decipher cancer genetic heterogeneity. *Nucleic Acids Res*. 2017;45(W1):W514-
457 W22.
- 458 42. Breuza L, Poux S, Estreicher A, Famiglietti ML, Magrane M, Tognolli M, et al. The UniProtKB
459 guide to the human proteome. *Database*. 2016;2016.
- 460 43. Burley SK, Bhikadiya C, Bi C, Bittrich S, Chen L, Crichlow GV, et al. RCSB Protein Data Bank:
461 powerful new tools for exploring 3D structures of biological macromolecules for basic and applied
462 research and education in fundamental biology, biomedicine, biotechnology, bioengineering and energy
463 sciences. *Nucleic acids research*. 2021;49(D1):D437-D51.
- 464 44. Chang MT, Asthana S, Gao SP, Lee BH, Chapman JS, Kandath C, et al. Identifying recurrent
465 mutations in cancer reveals widespread lineage diversity and mutational specificity. *Nat Biotechnol*.
466 2016;34(2):155-63.
- 467 45. Holm L, Laakso LM. Dali server update. *Nucleic acids research*. 2016;44(W1):W351-W5.
- 468 46. Forbes SA, Beare D, Boutselakis H, Bamford S, Bindal N, Tate J, et al. COSMIC: somatic cancer
469 genetics at high-resolution. *Nucleic acids research*. 2017;45(D1):D777-D83.
- 470 47. Landrum MJ, Lee JM, Benson M, Brown GR, Chao C, Chitipiralla S, et al. ClinVar: improving
471 access to variant interpretations and supporting evidence. *Nucleic acids research*. 2018;46(D1):D1062-
472 D7.

- 473 48. Griffith M, Spies NC, Krysiak K, McMichael JF, Coffman AC, Danos AM, et al. CIViC is a
474 community knowledgebase for expert crowdsourcing the clinical interpretation of variants in cancer.
475 *Nature genetics*. 2017;49(2):170-4.
- 476 49. Doncheva NT, Klein K, Domingues FS, Albrecht M. Analyzing and visualizing residue networks
477 of protein structures. *Trends in biochemical sciences*. 2011;36(4):179-82.
- 478 50. Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, et al. UCSF
479 Chimera—A visualization system for exploratory research and analysis. 2004;25(13):1605-12.
- 480 51. Holmås S, Riudavets Puig R, Acencio ML, Mironov V, Kuiper M. The Cytoscape BioGateway
481 App: explorative network building from an RDF store. Oxford University Press; 2020.
- 482 52. Stefl S, Nishi H, Petukh M, Panchenko AR, Alexov E. Molecular mechanisms of disease-causing
483 missense mutations. *Journal of molecular biology*. 2013;425(21):3919-36.
- 484 53. Capriotti E, Fariselli P, Casadio R. I-Mutant2. 0: predicting stability changes upon mutation from
485 the protein sequence or structure. *Nucleic acids research*. 2005;33(suppl_2):W306-W10.
- 486 54. Dehouck Y, Kwasigroch JM, Gilis D, Rooman M. PoPMuSiC 2.1: a web server for the
487 estimation of protein stability changes upon mutation and sequence optimality. *BMC bioinformatics*.
488 2011;12(1):1-12.
- 489 55. Sim N-L, Kumar P, Hu J, Henikoff S, Schneider G, Ng PC. SIFT web server: predicting effects
490 of amino acid substitutions on proteins. *Nucleic acids research*. 2012;40(W1):W452-W7.
- 491 56. Choi Y, Chan AP. PROVEAN web server: a tool to predict the functional effect of amino acid
492 substitutions and indels. *Bioinformatics*. 2015;31(16):2745-7.
- 493 57. Parthiban V, Gromiha MM, Schomburg D. CUPSAT: prediction of protein stability upon point
494 mutations. *Nucleic acids research*. 2006;34(suppl_2):W239-W42.
- 495 58. Pejaver V, Urresti J, Lugo-Martinez J, Pagel KA, Lin GN, Nam H-J, et al. Inferring the molecular
496 and phenotypic impact of amino acid variants with MutPred2. *Nature communications*. 2020;11(1):1-13.
- 497 59. Kodaz H, Kostek O, Hacıoglu MB, Erdogan B, Kodaz CE, Hacibekiroglu I, et al. Frequency of
498 RAS mutations (KRAS, NRAS, HRAS) in human solid cancer. *Breast cancer*. 2017;7(5).

- 499 60. Svensmark JH, Brakebusch C. Rho GTPases in cancer: friend or foe? *Oncogene*.
500 2019;38(50):7447-56.
- 501 61. Bertola D, Buscarilli M, Stabley DL, Baker L, Doyle D, Bartholomew DW, et al. Phenotypic
502 spectrum of Costello syndrome individuals harboring the rare HRAS mutation p. Gly13Asp. *American*
503 *Journal of Medical Genetics Part A*. 2017;173(5):1309-18.
- 504 62. Gripp KW, Baker L, Robbins KM, Stabley DL, Bellus GA, Kolbe V, et al. The novel duplication
505 HRAS c. 186_206dup p.(Glu62_Arg68dup): clinical and functional aspects. *European Journal of Human*
506 *Genetics*. 2020;28(11):1548-54.
- 507 63. Homami A, Kachoei ZA, Asgarie M, Ghazi F. Analysis of FGFR3 and HRAS genes in patients
508 with bladder cancer. *Medical Journal of the Islamic Republic of Iran*. 2020;34:108.
- 509 64. Pozdeyev N, Rose MM, Bowles DW, Schweppe RE, editors. *Molecular therapeutics for*
510 *anaplastic thyroid cancer*. *Seminars in cancer biology*; 2020: Elsevier.
- 511 65. Gamayunov BN, Korotkiy NG, Baranova EE. Phacomatosis pigmentokeratolica or the
512 Schimmelpenning-Feuerstein-Mims syndrome? *Clinical case reports*. 2016;4(6):564.
- 513 66. Kaur HB, Salles DC, Paulk A, Epstein JI, Eshleman JR, Lotan TL. PIN-like ductal carcinoma of
514 the prostate has frequent activating RAS/RAF mutations. *Histopathology*. 2021;78(2):327-33.
- 515 67. Witvliet DK, Strokach A, Giraldo-Forero AF, Teyra J, Colak R, Kim PM. ELASPIC web-server:
516 proteome-wide structure-based prediction of mutation effects on protein stability and binding affinity.
517 *Bioinformatics*. 2016;32(10):1589-91.
- 518 68. Papanikolaou N, Mantsou A, Kalosidis N. From driver mutations to driver cancer networks: Why
519 we need a new paradigm. *Cancer Studies*. 2018;2(1):1.
- 520 69. Scott F, Fala AM, Pennicott LE, Reuillon TD, Massirer KB, Elkins JM, et al. Development of 2-
521 (4-pyridyl)-benzimidazoles as PKN2 chemical tools to probe cancer. *Bioorganic & medicinal chemistry*
522 *letters*. 2020;30(8):127040.

- 523 70. Yang CS, Melhuish TA, Spencer A, Ni L, Hao Y, Jividen K, et al. The protein kinase C super-
524 family member PKN is regulated by mTOR and influences differentiation during prostate cancer
525 progression. *The Prostate*. 2017;77(15):1452-67.
- 526 71. Muñoz-Maldonado C, Zimmer Y, Medová M. A comparative analysis of individual RAS
527 mutations in cancer biology. *Frontiers in oncology*. 2019;9:1088.
- 528 72. Schaefer A, Reinhard NR, Hordijk PL. Toward understanding RhoGTPase specificity: structure,
529 function and local activation. *Small GTPases*. 2014;5(2):e968004.
- 530 73. Mosteller RD, Han J, Broek D. Identification of residues of the H-ras protein critical for
531 functional interaction with guanine nucleotide exchange factors. *Molecular and Cellular Biology*.
532 1994;14(2):1104-12.
- 533 74. Kakiuchi M, Nishizawa T, Ueda H, Gotoh K, Tanaka A, Hayashi A, et al. Recurrent gain-of-
534 function mutations of RHOA in diffuse-type gastric carcinoma. *Nature genetics*. 2014;46(6):583-7.
- 535 75. Aoki Y, Niihori T, Kawame H, Kurosawa K, Ohashi H, Tanaka Y, et al. Germline mutations in
536 HRAS proto-oncogene cause Costello syndrome. *Nature genetics*. 2005;37(10):1038-40.
- 537 76. Sherry ST, Ward M-H, Kholodov M, Baker J, Phan L, Smigielski EM, et al. dbSNP: the NCBI
538 database of genetic variation. *Nucleic acids research*. 2001;29(1):308-11.
- 539